

Test Data Likelihood for PLSA Models

Thorsten Brants
Palo Alto Research Center (PARC)
3333 Coyote Hill Rd, Palo Alto, CA 94304, USA
brants@parc.com

Abstract

Probabilistic Latent Semantic Analysis (PLSA) is a statistical latent class model that has recently received considerable attention. In its usual formulation it cannot assign likelihoods to unseen documents. Furthermore, it assigns a probability of zero to unseen documents during training. We point out that one of the two existing alternative formulations of the Expectation-Maximization algorithms for PLSA does not require this assumption. However, even that formulation does not allow calculation of the actual likelihood values. We therefore derive a new test-data likelihood substitute for PLSA and compare it to two existing likelihood substitutes. An empirical evaluation shows that our new likelihood substitute produces the best predictions about accuracies in two different IR tasks and is therefore best suited to determine the number of EM step when training PLSA models. The new likelihood measure and its evaluation also suggest that PLSA is not very sensitive to overfitting for the two tasks considered. This renders additions like tempered EM that especially address overfitting unnecessary.

Keywords: Probabilistic Latent Semantic Analysis, Likelihood

1 Introduction

Probabilistic Latent Semantic Analysis is a statistical latent class model (or aspect model) that recently has shown excellent results in several IR related tasks [5, 10, 8]. It has the positive feature of assigning probability distributions over classes to documents and words, unlike other clustering techniques that do a hard assignment to one class. Aspect models are also successfully used in other areas of natural language processing, like language modeling [13] or parsing [12].

However, PLSA does not provide a way of assigning likelihoods to instances of the observed variables that do not occur in the training corpus. In the case of retrieval tasks, it cannot assign likelihoods to unseen test documents. Even less desirable, its usual formulation proposes zero probabilities for unseen documents during training and then ignores this fact during testing when calculating probabilities conditioned on the test document.

As we will show, one of the two alternative formulations of the Expectation-Maximization algorithm for PLSA does not assume zero-probabilities for test documents, which favors this version. But still there is no intrinsic way of likelihood calculations. We therefore derive a new test data likelihood substitute and compare it empirically with two existing and previously used likelihood substitutes. The comparison reveals a good qualitative match between the new likelihood substitute and task results, which makes it best suited to determine the number of EM steps for training PLSA in an unsupervised setting.

As an additional result of our comparison we find evidence that PLSA is very robust against overfitting for the two tasks considered. On the one hand, aspect models have previously been reported to be very sensitive to overfitting when applied to information retrieval tasks [2, 5, 10, 8] even with small numbers of classes. Therefore, the studies usually proposed using tempered EM in order to avoid overtraining. On the other hand, studies using aspect models for language modeling and parsing report these models to be relatively robust against overtraining when the number of classes is much smaller than the size of the vocabulary [13, 12]. Our experiment plotting the number of EM steps vs. retrieval and segmentation task results show that PLSA is not sensitive to overfitting in these two tasks and therefore render techniques like tempered EM for PLSA unnecessary.

The rest of the paper is structured as follows. Section 2 shortly introduces the PLSA model. Section 3.2 presents an alternative formulation and points out its advantage of avoiding zero probabilities for test documents. Section 4 introduces our new method of calculating substitutes for PLSA test documents likelihoods and also presents two existing methods. Section 5 empirically compares the three likelihood measures against each other and against results in document retrieval and text segmentation. Section 6 contains conclusions.

2 The PLSA Model

Probabilistic Latent Semantic Analysis (PLSA) is a statistical latent class model or aspect model [8, 7]. The model is fitted to a training corpus by the Expectation Maximization (EM) algorithm [4]. It assigns probability distributions over classes to words and documents and thereby allows them to belong to more than one class, and not to only one class as is true of most other classification methods. PLSA represents the joint probability of a document d and a word w based on a latent class variable z :¹

$$P(d, w) = P(d) \sum_z P(w|z)P(z|d) \quad (1)$$

PLSA has the following view of how a document is generated: first a document $d \in \mathcal{D}$ (i.e., its dummy label) is chosen with probability $P(d)$. For each word in document d , a latent topic $z \in \mathcal{Z}$ is chosen with probability $P(z|d)$, which in turn is used to choose the word $w \in \mathcal{W}$ with probability $P(w|z)$.

A model is fitted to a document collection \mathcal{D} by maximizing the log-likelihood function \mathcal{L} :

$$\mathcal{L} = \sum_{d \in \mathcal{D}} \sum_{w \in d} f(d, w) \log P(d, w). \quad (2)$$

The E-step in the EM-algorithm is

$$P(z|d, w) = \frac{P(z)P(d|z)P(w|z)}{\sum_{z'} P(z')P(d|z')P(w|z')} \quad (3)$$

and the M-step consists of

$$P(w|z) = \frac{\sum_d f(d, w)P(z|d, w)}{\sum_{d, w'} f(d, w')P(z|d, w')} \quad (4)$$

$$P(d|z) = \frac{\sum_w f(d, w)P(z|d, w)}{\sum_{d', w} f(d', w)P(z|d', w)} \quad (5)$$

$$P(z) = \frac{\sum_{d, w} f(d, w)P(z|d, w)}{\sum_{d, w} f(d, w)}. \quad (6)$$

The parameters are either randomly initialized or according to some prior knowledge.

¹Unless otherwise noted, we use the following notational conventions: training documents $d, d' \in \mathcal{D}$, test documents $q, q' \in \mathcal{Q}$, words $w, w' \in \mathcal{W}$, and classes $z, z' \in \mathcal{Z}$.

The parameters $P(w|z)$ obtained in the training process are used to calculate $P(z|q)$ for new documents q (queries) with the folding-in process. Folding-in uses Expectation-Maximization as in the training process; the E-step is identical, the M-step keeps all the $P(w|z)$ constant and re-calculates $P(z|q)$. Usually, a very small number of iterations is sufficient for folding-in.

3 Alternative PLSA Formulation

3.1 Zero-Probability Documents

Note that the EM process in (3) – (6) is based on the probabilities $P(d|z)$ while the probabilities used for folding-in of query documents are $P(z|q)$. Bayes’ formula allows the reverse conditioning

$$P(z|q) = \frac{P(z)P(q|z)}{P(q)}. \quad (7)$$

A subtle but crucial point of PLSA is that we do not know $P(q)$. Even worse, the model as it is usually formulated [8] (and as it is shortly presented in the previous section) postulates $P(q|z) = 0$ for all $z \in \mathcal{Z}$, and hence $P(q) = 0$, for all $q \notin \mathcal{D}$, where \mathcal{D} is the training collection. This can be derived from equation 5. For any given $z \in \mathcal{Z}$, we have

$$\sum_{d \in \mathcal{D}} P(d|z) = \frac{\sum_{w \in \mathcal{W}} f(d, w)P(z|d, w)}{\sum_{d' \in \mathcal{D}} \sum_{w \in \mathcal{W}} f(d', w)P(z|d', w)} \quad (8)$$

$$= 1. \quad (9)$$

Since all the probability mass is on the $d \in \mathcal{D}$, it follows that $P(q|z) = 0$ for any $q \notin \mathcal{D}$.

The PLSA formulation based on $P(d|z)$ simply ignores this zero-probability problem and proceeds with the folding-in process. It calculates probabilities for $P(z|q)$ during folding-in and avoids looking at $P(q)$, hence it ignores the division of 0 by 0. The deviation by zero is not made explicitly during folding-in, but it is there implicitly because of equation 7.

3.2 PLSA Reformulation

We now switch to a formulation of EM for PLSA that was used in [5]. That paper presented the reformulation without giving a motivation for preferring one version or the other, and without addressing implications.

We now use probabilities $P(z|d)$ instead of $P(d|z)$. The E-Step is

$$P(z|d, w) = \frac{P(z|d)P(w|z)}{\sum_{z'} P(z'|d)P(w|z')} \quad (10)$$

and the M-step consists of

$$P(w|z) = \frac{\sum_d f(d, w)P(z|d, w)}{\sum_{d, w'} f(d, w')P(z|d, w')} \quad (11)$$

$$P(z|d) = \frac{\sum_w f(d, w)P(z|d, w)}{\sum_{w, z'} f(d, w)P(z'|d, w)} \quad (12)$$

$$P(z) = \frac{\sum_{d, w} f(d, w)P(z|d, w)}{\sum_{d, w} f(d, w)} \quad (13)$$

This section investigates the differences between the two formulations. The log-likelihood function (2) can be written as

$$\mathcal{L} = \sum_{d \in \mathcal{D}} \sum_{w \in d} f(d, w) \log P(d) + \sum_{d \in \mathcal{D}} \sum_{w \in d} f(d, w) \log P(w|d) \quad (14)$$

because of

$$\log P(d, w) = \log P(d)P(w|d) = \log P(d) + \log P(w|d). \quad (15)$$

Both parts of equation (14) can be optimized independently of each other (cf. [10]). The PLSA parameters $P(w|z)$ and $P(z|d)$ only occur in the right hand part (in $P(w|d)$). Therefore, it is sufficient to maximize that part, which is done by equations (10) to (12).

EM for PLSA based on $P(z|d)$ does not use probabilities $P(d)$, and hence $P(q)$, at all. They are not part of equations (10) to (12). The formulas only require $P(d) > 0$ for all $d \in \mathcal{D}$ in order to obtain valid probabilities conditioned on d . And they require

$$\sum_{d \in \mathcal{D}} P(d) < 1 \quad (16)$$

in order to reserve probability mass to unseen documents q (i.e., having $P(q) > 0$ for test documents q). We still do not have a way of calculating the probability value for unseen documents, but the model now accounts for them and does no longer postulate zero probabilities during training. Knowledge of the exact value of $P(d)$ is not necessary during training because it does not show up in the EM steps. We only need probabilities conditioned on d , so we assume $P(d) > 0$ and $P(q) > 0$.

With this formulation, folding-in uses the same EM formulas as training. For a new document q , the algorithm alternates between the E-step for $P(z|q, w)$ (10) and the M-step for $P(z|q)$ (12). Note that both $P(z|q, w)$ and $P(z|q)$ are independent of any $q' \neq q$, i.e., folding-in of q proceeds without information about other test documents q' .

To sum up, the reformulation maximizes the conditional part of the likelihood function in equation (14). Therefore, it does not need to refer to $P(d)$ and can assume $P(q) > 0$ for documents q not in the training set. To our knowledge, this motivation for the particular PLSA model formulation has not been published before.

This is only half of the solution to the PLSA likelihood problem. We know that $P(q) > 0$, but we still do not know the actual value of $P(q)$. However, knowledge of non-training data likelihoods is very useful when determining the number of EM-steps on a held-out data set during unsupervised learning.

4 Likelihood Calculation

The EM formulation for PLSA in section 3.2 eliminates the weakness of assigning zero probabilities to unseen documents. However, it does so by avoiding to assign a particular likelihood to any of the documents. Neither the usual formulation nor the reformulation can assign a likelihood to an unseen query document. Previously, two different test data likelihood substitutes have been suggested, of which one has been used in previous publications. In this section, we first present a new (third) method to calculate likelihood substitutes for unseen documents. This method is derived from the PLSA likelihood maximization definition. The method does not give us the true likelihood values, but we show that it is maximized at the same points in the parameter space as the true likelihood. We then present the two existing likelihood substitutes.

For the sake of brevity, we make two simplifying assumptions without losing generality. First, we do not use logarithms here (for an implementation, however, logarithmic values need to be used). Second, we assume each w occurs exactly once in a document q . If a term occurs several times, it is simply repeated in q . This saves us from keeping track of frequency counts $f(q, w)$ and lets us concentrate on the differences of the three alternatives.

4.1 Likelihood Based on Folding-In

The overall goal is to calculate a likelihood $p_{test}(\mathcal{Q}|\theta)$ for a test collection \mathcal{Q} conditioned on the model parameters θ , which is in analogy to (2):

$$p_{test}(\mathcal{Q}|\theta) = \prod_{q \in \mathcal{Q}} \prod_{w \in q} P(q, w|\theta). \quad (17)$$

In the following, we will omit θ from the notation. (17) can be re-written as

$$p_{test}(\mathcal{Q}) = \prod_{q \in \mathcal{Q}} \prod_{w \in q} P(q)P(w|q) \quad (18)$$

$$= \prod_{q \in \mathcal{Q}} \left[P(q)^{|q|} \prod_{w \in q} P(w|q) \right] \quad (19)$$

where $|q|$ is the number of words in test document q . The difficult part is coming up with a value for $P(q)$. However, the following re-formulation allows us to eliminate $P(q)$ from the calculation.

Further separating the elements in (19) yields

$$p_{test}(\mathcal{Q}) = \left[\prod_{q \in \mathcal{Q}} P(q)^{|q|} \right] \prod_{q \in \mathcal{Q}} \prod_{w \in q} P(w|q) \quad (20)$$

$$= \left[\prod_{q \in \mathcal{Q}} P(q)^{|q|} \right] \prod_{q \in \mathcal{Q}} \prod_{w \in q} \left[\sum_z P(w|z)P(z|q) \right] \quad (21)$$

The left part in square brackets is not dependent on the model parameters $P(w|z)$ and $P(z|d)$, and also not on $P(z|q)$. Therefore, the left part and the rest can be maximized independently in order to find the product's maximum. In case we are interested in finding the model parameters that yield the maximization it is sufficient to calculate the right (conditional) part

$$p_{test}^{cond}(\mathcal{Q}) = \prod_{q \in \mathcal{Q}} \prod_{w \in q} \left[\sum_z P(w|z)P(z|q) \right] \quad (22)$$

Probabilities $P(w|z)$ are the model parameters, $P(z|q)$ for test documents q can be derived by folding-in. Our proposed method of generating a likelihood substitute $p_{test}^f(\mathcal{Q})$ for an unseen test collection \mathcal{Q} uses

(22):

$$p_{test}^{fi}(\mathcal{Q}) = \prod_{q \in \mathcal{Q}} \prod_{w \in q} \sum_z P(w|z) P_{fi}(z|q) \quad (23)$$

$$= \prod_{q \in \mathcal{Q}} \prod_{w \in q} P_{fi}(w|q) \quad (24)$$

P_{fi} marks the probabilities from folding-in, the others are obtained from the training process. p_{test}^{fi} does not represent the true likelihood value, but if our goal is likelihood maximization we find exactly the same parameter setting as if we had maximized the true likelihood, which follows from the decomposition in (21).

4.2 Likelihood Based on Marginalizing over the Training Documents

Another method estimates $P(q)$ by marginalizing over all training documents d , yielding:

$$p_{test}^d(\mathcal{Q}) = \prod_{q \in \mathcal{Q}} \sum_{d \in \mathcal{D}} P(d) \prod_{w \in q} \sum_z P(w|z) P(z|d) \quad (25)$$

$$= \prod_{q \in \mathcal{Q}} \sum_{d \in \mathcal{D}} P(d) \prod_{w \in q} P(w|d) \quad (26)$$

This method was used in [2] for the comparison of different models. Its view is that a test document q is an “average training document”, i.e., it calculates the probability that a particular training document d generates all words w of the test document and then averages over all training documents. Likelihood based on marginalizing over the training documents conditions word probabilities on each training document separately and then averages. Unfortunately, this does not match very well the view of PLSA that a document is generated by its own (and possibly unique) mixture of topics z . We therefore expect this likelihood estimate in general to be too low.

This method of calculating a likelihood substitute was most probably also used in the work by Hofmann *et al.* Their papers do not directly mention the formula, but their plots of the perplexity (which is directly related to the log likelihood) are equivalent to those that we will present in section 5, e.g. those presented in Figures 5 and 6 of [10], or those presented in Figure 7 of [9].

4.3 Likelihood Based on a Word Unigram Model

The third method, mentioned in [2], creates a word unigram model by marginalizing out d and z from the PLSA model and then calculates test document likelihood by multiplying the word probabilities:

$$p_{test}^w(\mathcal{Q}) = \prod_{q \in \mathcal{Q}} \prod_{w \in q} \sum_{d \in \mathcal{D}} P(d) \sum_z P(w|z)P(z|d) \quad (27)$$

$$= \prod_{q \in \mathcal{Q}} \prod_{w \in q} P(w) \quad (28)$$

As with the previous methods, this way of calculating a likelihood ignores the PLSA view of how a document is created: select a document, select a mixture of topics for the document, select the words based on the mixture of topics.

4.4 Concluding Remarks on the Three Methods

PLSA does not provide a way of calculating test data likelihood. Our newly presented method circumvents the problem by calculating only a part of the likelihood, but this part is maximized by the same arguments that maximize the true likelihood. The other two try to circumvent the problem by marginalizing over one or two parameters. Likelihood based on folding-in follows the view of PLSA of how a document is generated ($d \rightarrow z \rightarrow w$). The other two methods deviate substantially from the PLSA generative view and we therefore expect the folding-in likelihood to be best suited.

Sections 4.1 to 4.3 presented the likelihood methods using notational simplifications. Practical applications would perform the calculations using logarithms, and they would not simply repeat multiple occurrences of the same word but directly model the frequency $f(d, w)$. We therefore repeat the three likelihood formulas, but now are using logarithms and $f(d, w)$.

Folding-in log-likelihood (cf. 24):

$$\mathcal{L}_{test}^{fi}(\mathcal{Q}) = \sum_{q \in \mathcal{Q}} \sum_{w \in q} f(q, w) \log \sum_z P(w|z)P_{fi}(z|q) \quad (29)$$

Document-marginal log-likelihood (cf. 26):

$$\mathcal{L}_{test}^d(\mathcal{Q}) = \sum_{q \in \mathcal{Q}} \log \sum_{d \in \mathcal{D}} P(d) \prod_{w \in q} \left[\sum_z P(w|z)P(z|d) \right]^{f(q,w)} \quad (30)$$

Document-marginal log-likelihood is the most difficult to calculate of the three methods. Because of the sum within the logarithm the product over all words in q needs to be actually calculated, which means one has to take care of floating point underflows for long test documents.

Word unigram log-likelihood (cf. 28):

$$\mathcal{L}_{test}^w(\mathcal{Q}) = \sum_{q \in \mathcal{Q}} \sum_{w \in q} f(q, w) \log \sum_{d \in \mathcal{D}} P(d) \sum_z P(w|z)P(z|d) \quad (31)$$

5 Experiments

The experiments reported in this section aim at aligning the three likelihood measures presented in the previous section with error rates obtained in information retrieval tasks. The first task is a standard retrieval task, the second task evaluates on text segmentation.

5.1 Corpora and Pre-Processing

For the retrieval task, we used the MED collection, consisting of 1,033 document abstracts from the National Library of Medicine, 15 queries, and manually assigned relevance judgements for each of the 15 queries.

For the segmentation task, we used the ModApte split of the Reuters-21578 Corpus. The training set consists of 7769 news articles with approx. 1.1 million tokens. The test set consists of 3019 articles with approx. 400,000 tokens. The corpus was prepared according to the description in [11]. Exactly two original documents from the test part are concatenated to form one test document, one from one of the 10 larger categories, one from a randomly chosen other category. In total, 500 test documents are generated. The task is to find the boundary of the two original documents. More details about the segmentation task can be found in [6, 3]. The error metric is described in [1].

In both cases, the corpora were stemmed, stopwords and words with frequency 1 were removed, the rest was mapped to lower case.

5.2 Likelihood Calculation

We performed all experiments with each of the two corpora using numbers of classes ranging from 16 to 8,192 in powers of 2, and repeated all experiments using several different random EM initializations. Results using the different numbers of classes and initializations were very similar. We therefore present as representatives details of the experiments

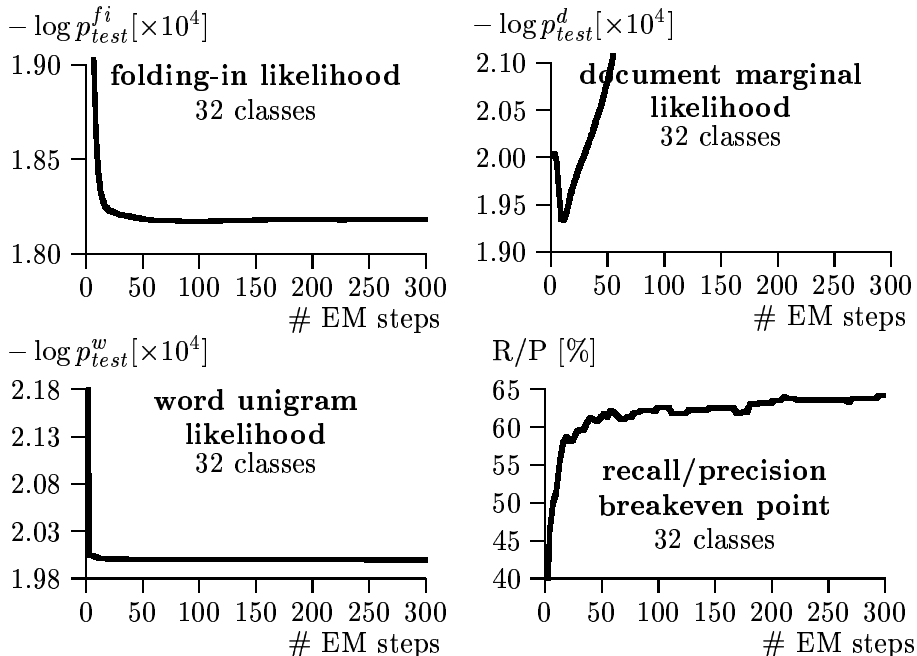


Figure 1: MED Corpus Likelihood and Retrieval Performances, 32 classes

with 32 classes and 128 classes. The PLSA models were randomly initialized. After each EM step, we recorded the three likelihood measures for the test set as presented in section 4. Figures 1 to 4 show the negative log likelihoods for steps 1 to step 300 (please note that lower values mean higher likelihood and vice versa since we are plotting the negative log likelihood).

The folding-in log-likelihood \mathcal{L}_{test}^{fi} slightly decreases after an initial sharp increase for the MED corpus. For the Reuters Corpus, it monotonically increases even after a large number of iterations. The difference may be attributed to the different corpus sizes. The Reuters Corpus is approx. seven times larger than the MED corpus, making overfitting less likely for the Reuters corpus.

The document marginal log-likelihood \mathcal{L}_{test}^d sharply increases for approx. 10 initial iterations, and then rapidly decreases towards 0. (with the exception of the graph for the Reuters Corpus with 32 classes is slightly different; we continued the experiment and found that likelihood in this case will finally also decrease below random level). In the other cases, the likelihood is down to random level after 20 – 80 EM iterations, i.e. at the same level where it was after initialization. The steep

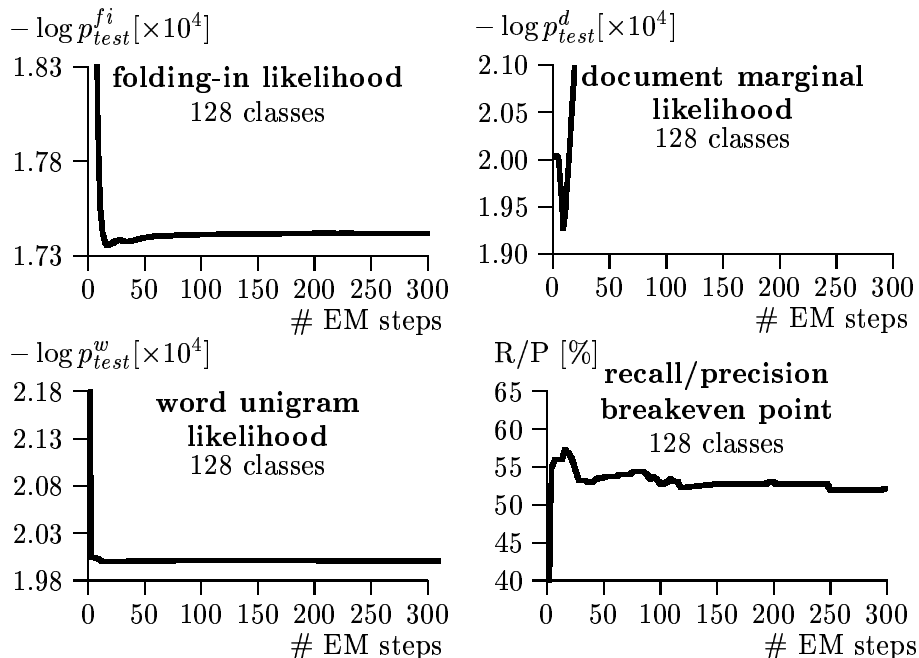


Figure 2: MED Corpus Likelihoods and Retrieval Performance, 128 classes

decrease in document marginal likelihood can be explained by equation 26 and that with larger numbers of classes PLSA tends to move most of the probabilities $P(w|z)$ and $P(z|d)$ very close to zero². This will drop a large number of $P(w|d)$ to zero and finally assign a probability of 0 to all test documents unless they are (almost) identical to one of the training documents. The flaw is that this likelihood uses the probability that any training document has generated the vector of words of the test document $P(w|d)$, but what we are really interested in is the probability that the test document has generated them $P(w|q)$.

Word unigram likelihood \mathcal{L}_{test}^w increases sharply from random level for only very few EM iterations (in some cases for one iteration only). The Reuters Corpus graph shows a small decrease around 10 iterations, the MED graph levels out immediately. Using this as a model performance measure would mean that the model is overtrained already after just one step in the Reuters case, and after very few steps in the MED case. This way of calculating likelihood suffers from its very restricted underlying model, word unigrams, which is not the way PLSA assigns

²“Very close to zero” becomes zero given the finite precision of any computation.

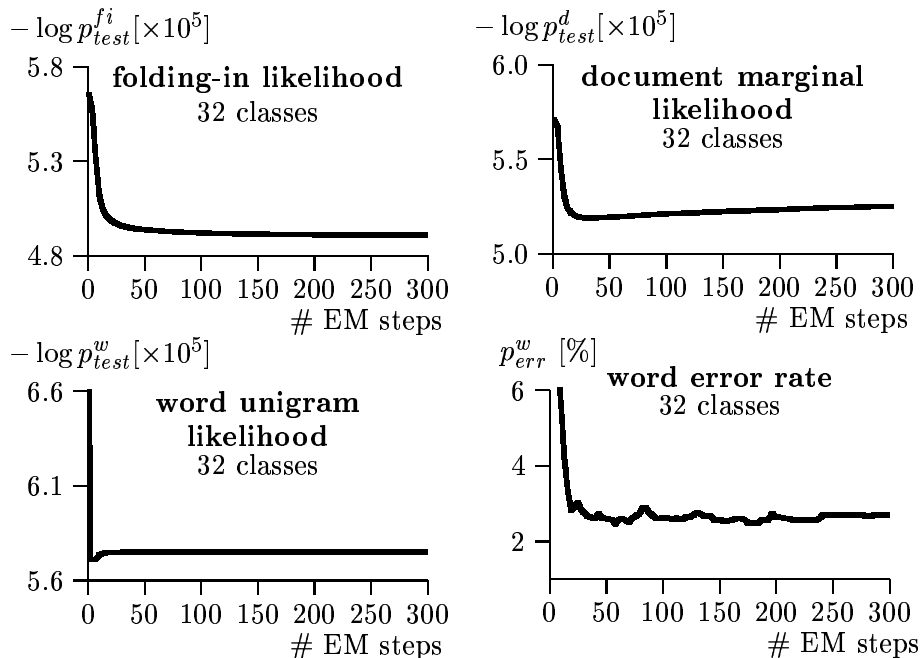


Figure 3: Reuters-21578 Corpus Likelihoods and Word Error Rate, 32 classes

probabilities.

5.3 Recall/Precision and Error Rate

We now present results on IR tasks for PLSA models that are trained with varying numbers of EM iterations. The goal is a qualitative comparison with the three methods of likelihood calculation. Therefore, we do not use additional methods like tempered EM or interpolation with other models to improve results but report plain PLSA results.

The task performed on the MED corpus is document retrieval. The test set consists of 15 queries, each consisting of one line of text. The goal is to find relevant documents in the collection. Automatically retrieved documents are compared against manual judgements on the documents. We report the recall/precision breakeven point, i.e., we take as many documents from the top of the ranked retrieved list of documents until recall and precision have the same value. The recall/precision breakeven point is plotted against the number of EM used for training the PLSA model. Results are shown at the bottom of figures 1 and 2.

We do not find any overfitting during the first 300 iterations for the

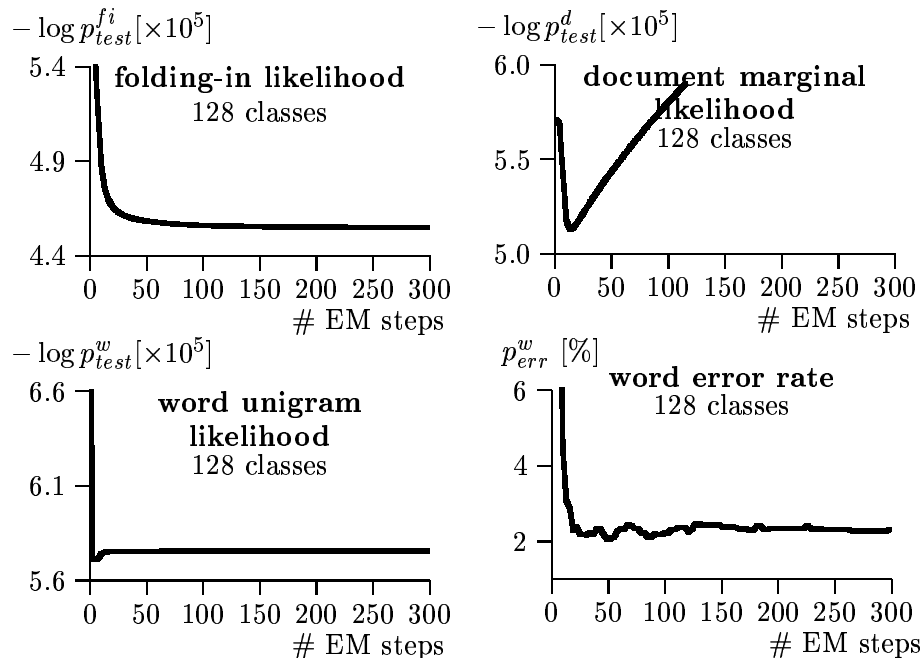


Figure 4: Reuters-21578 Corpus Likelihoods and Word Error Rate, 128 classes

model with 32 classes. For 128 classes, the model improves until around 20 iterations, then slightly drops and stays relatively constant after 50 iterations. This result is a surprise since PLSA was previously reported to require tempered EM to avoid overfitting (e.g., [8, 2]).

The second task is text segmentation, performed on the Reuters corpus. Results (error rates) are shown in the bottom part of figures 3 and 4. Again, we find improving results until around 20 iterations. After this point word error rate stays relative constant. There is no sign of overtraining within 300 iterations.

Qualitatively comparing the graphs for the three likelihood measures with the retrieval and segmentation results, the curves of folding-in likelihood and task performance match best. Document marginal likelihood and word unigram likelihood have very different curves. Folding-in likelihood even predicts the slight decrease in recall/precision around iteration 20 for the retrieval task with 128 classes. The newly presented folding-in likelihood was derived from the PLSA likelihood definition as it is used for the Expectation-Maximization algorithm and empirically matches recall/precision and error rates very well, while the other two

are ad-hoc substitutes that do not share the generative view of PLSA.

Our empirical results might also hint towards an explanation why Hofmann *et al.* found better results for tempering than for early stopping in their experiments on the MED corpus. If they actually used the document marginal likelihood, maximum likelihood on the test or held-out data was achieved too early. As an example, in figure 2, document marginal log-likelihood reaches the maximum of -19,270 after 9 iterations. At this point, the recall-precision break-even point is at 56.0%. Recall/precision maximum is at a later point: it is 57.2% after 16 iterations. When using the document-marginal likelihood, tempering seems necessary in order to go beyond the too early stopping point. However, folding-in likelihood does not need this heuristic addition of tempering. It achieves the maximum of -17,358 after 17 iterations and therefore comes very close to the optimal performance. Results similarly favor folding-in likelihood in the other cases: document marginal likelihood has its maximum too early in most of the cases. Tempering might still be useful in order to escape unfavorable local maxima, but it is not necessary to avoid overfitting in the two presented tasks.

6 Conclusions

Investigating the two alternative formulations of EM for PLSA we observed that the use of $P(z|d)$ instead of $P(d|z)$ removes the weakness of assuming zero probabilities for unseen documents. However, both formulations do not assign a particular likelihood value to test documents which makes it necessary to resort to substitutes.

We presented a new test data likelihood substitute that we derived from the PLSA likelihood definition as used for the Expectation-Maximization algorithm. This was named the folding-in likelihood.

We plotted graphs of the folding-in likelihood as well as two other likelihood substitutes: the document-marginal likelihood and the word unigram likelihood. Folding-in likelihood suggests that there is very little overfitting for the two corpora considered, document marginal likelihood suggests that there is heavy overfitting, word unigram likelihood seems unusable because it drops for very few (sometimes only one) iterations and then stays almost constant.

Comparison to recall/precision in a retrieval task and error rate in a segmentation task showed that the best predictions are made by folding-in likelihood. This suggest that folding-in likelihood is well suited to determine the number of EM steps in an unsupervised learning setting.

Claims in the literature that PLSA models in retrieval are very likely

to overfit when using (untempered) EM cannot be confirmed; more to the contrary, PLSA models seem very robust. To our knowledge, this is the first investigation that systematically increases the number of EM steps for a PLSA model and directly reports recall/precision/error rate after each step. Previous publications used the test set perplexity instead (which is directly related to the likelihood), but, as we showed, two of the three presented likelihood measures are poor performance predictors.

References

- [1] D. Beeferman, A. Berger, and J. Lafferty. Statistical models for text segmentation. *Machine Learning*, 34:177, 1999.
- [2] D. Blei, A. Ng, and M. Jordan. Latent dirichlet allocation. In *Proceedings of NIPS-2001*, Vancouver, BC, Canada, 2001.
- [3] F. Y. Y. Choi, P. Wiemer-Hastings, and J. Moore. Latent semantic analysis for text segmentation. In L. Lee and D. Harman, editors, *Proceedings of the 2001 Conference on Empirical Methods in Natural Language Processing*, pages 109–117, 2001.
- [4] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, 39(1):1–21, 1977.
- [5] D. Gildea and T. Hofmann. Topic based language models using em. In *Proceedings of 6th European Conference On Speech Communication and Technology (Eurospeech'99)*, Budapest, Hungary, 1999.
- [6] M. A. Hearst. Texttiling: Segmenting text into multi-paragraph subtopic passages. *Computational Linguistics*, 23(1):33–64, 1997.
- [7] T. Hofmann. Probabilistic latent semantic analysis. In *Proceedings of Uncertainty in Artificial Intelligence*, Stockholm, Sweden, 1999.
- [8] T. Hofmann. Probabilistic latent semantic indexing. In *Proceedings of SIGIR-99*, pages 35–44, Berkeley, CA, 2000.
- [9] T. Hofmann. Unsupervised learning by probabilistic latent semantic analysis. *Machine Learning*, 42:177–196, 2001.
- [10] T. Hofmann and J. Puzicha. Unsupervised learning from dyadic data. Technical Report TR-98-042, ICSI, Berkeley, CA, 1998.
- [11] H. Li and K. Yamanishi. Topic analysis using a finite mixture model. In *Proceedings of Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, pages 35–44, 2000.
- [12] M. Rooth, S. Riezler, D. Prescher, G. Carroll, and F. Beil. Inducing a semantically annotated lexicon via EM-based clustering. In *Proceedings of ACL-99*, College Park, MD, USA, 1999.
- [13] L. Saul and F. Pereira. Aggregate and mixed-order Markov models for statistical language processing. In *Proceedings of the Second Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 81–89. Association for Computational Linguistics, San Francisco, CA, 1997.