

4F13 - COURSEWORK #3: LATENT DIRICHLET ALLOCATION

CCN: 5654F

1. TASK A: MULTINOMIAL MAXIMUM LIKELIHOOD

The maximum likelihood multinomial assumes each word is drawn from one categorical distribution: $w_{nd} \sim \text{Cat}(\beta)$, where β is a vector containing the probability each word will be used. The maximum likelihood (ML) estimate for β for training set \mathcal{A} is simply the empirical word frequency for each word, as seen in Equation 1, where \mathbf{c} is the vector of counts of distinct word occurrences in the document set. \mathbf{c} is calculated as shown in Listing 1, and then divided by the total word count over the document, giving us $\hat{\beta}_{ML}$.

$$(1) \quad \begin{aligned} \hat{\beta}_{ML} &= \underset{\beta}{\operatorname{argmax}} \log \prod_{d=1}^D \prod_n^{N_d} \text{Cat}(w_{nd} | \beta) \\ &= \underset{\beta}{\operatorname{argmax}} \log \text{Mult}(c_1, \dots, c_M | \beta, N) = \frac{\mathbf{c}}{N} \end{aligned}$$

```
1 word_counts = np.zeros(W, dtype=int)
2 for doc_num in range(D):
3     doc_word_indices = np.where(A[:, 0] == doc_num + 1)
4     word_indices = np.array(A[doc_word_indices, 1])
5     word_counts[word_indices - 1] += np.array(A[doc_word_indices, 2])
```

LISTING 1. Task A code excerpts

The most and least probable words are then plotted in two bar plots, as shown in Figure 1. The highest log probability possible is generated when the test set contains only the most probable word: “bush”. For a test set L words long this is a log probability of $\log \hat{\beta}_{ML, \max}^L = L \log 0.0141 \approx -4.3L$ (0.0141 is the probability of “bush”). The lowest log probability occurs when any of the words in the test set did not appear in the training set, giving a probability of 0 and a log probability of $-\infty$. This is a very oversimplistic model as it cannot handle words that do not appear in the training set.

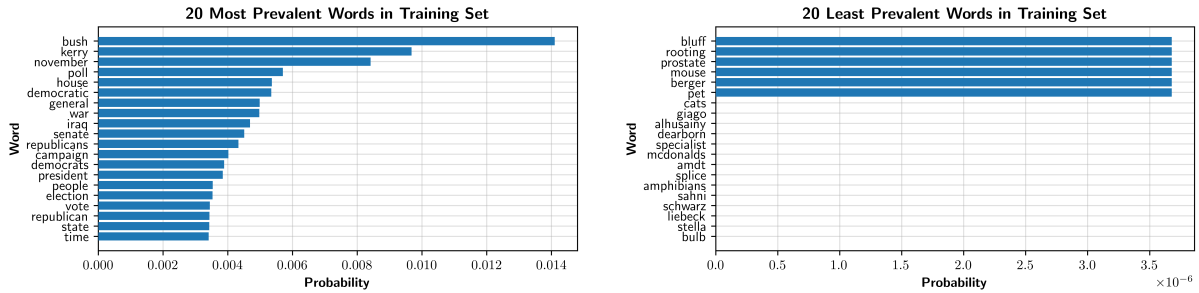


FIGURE 1. Task A Figures

2. TASK B: DIRICHLET PRIOR

The model was improved by adding a symmetric Dirichlet Prior: $\beta \sim \text{Dir}(\beta, \alpha)$, with concentration parameter α . Using Bayes Rule, the posterior distribution is calculated as in Equation 2, and as this is also a Dirichlet (due to the Dirichlet being a conjugate prior to the categorical), the maximum a posteriori (MAP) predictive distribution is as in Equation 3. The effect of the prior is to add a pseudocount of α to each count, by adding the lines in Listing 2.

$$(2) \quad P(\beta|\mathcal{A}) \propto P(\beta) \cdot P(\mathcal{A}|\beta) \propto \prod_{m=1}^M \beta_M^{c_m} \prod_{m=1}^M \beta_M^{c_m + \alpha - 1}$$

$$= \text{Dir}(\beta | \mathbf{c} + \alpha \cdot \mathbf{1})$$

$$(3) \quad \hat{\beta}_{MAP} = \frac{\mathbf{c} + \alpha \cdot \mathbf{1}}{\sum_{l=1}^M c_l + \alpha}$$

```

1 prior_counts = alpha * np.ones(W)
2 posterior_counts = word_counts + prior_counts
3 total_count = np.sum(posterior_counts)

```

LISTING 2. Task B code excerpts

Figure 2 shows different posterior predictive distributions for different values of α , for the top 20 and bottom 20 words (for $\alpha = 10$). The change is clear for the bottom 20 words: unseen words in training have a non-zero probability of occurring. As α increases the word probabilities become more uniform across the vocabulary. When $\alpha = 0$ the prior does nothing.

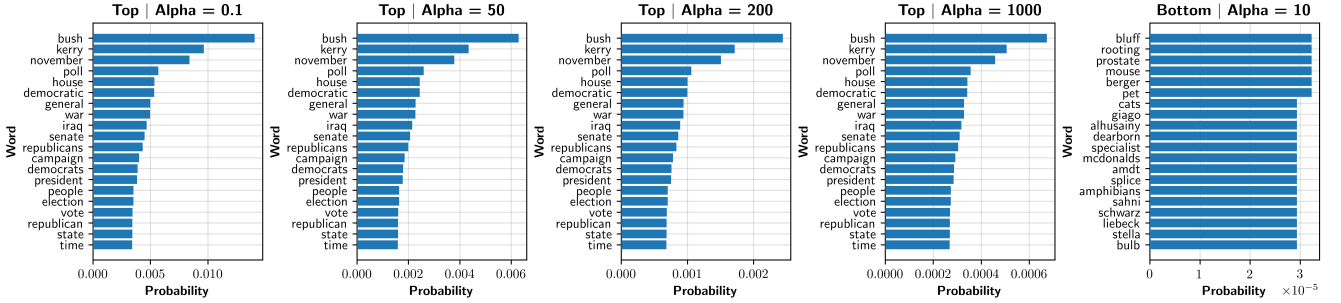


FIGURE 2. Dirichlet Prior - Posterior Predictive Distributions

3. TASK C: LOG PROBABILITIES

To calculate the log probability $L(d)$ of a test document d using the Bayesian Model, Equation 4 is used. $\{w_{nd}^*\}_{n=1}^{N_d}$ is the set of words in d . c_m^* is the count of word m in d . The data we are given however is just the word counts for each document, not the word order, which is obviously an important part of document analysis. To account for this a combinatorial term could be added as in Equation 5. However, we can consider all word orderings of documents to be equally likely for the moment and therefore this can be ignored and the categorical distribution function is used.

With α set to 0.1, the log probability of Document 2001 was calculated to be -3691.2, with 440 total words.

Also of use is the per-word perplexity, which is a measure of performance, and the definition is given in Equation 6. Listing 3 shows how the log probability and perplexity for Document 2001 were calculated.

$$(4) \quad L(d) = \log P(\{w_{nd}^*\}_{n=1}^{N_d} | \mathcal{A}) = \log \prod_{n=1}^{N_d} P(w_{nd}^* | \mathcal{A}) = \sum_{n=1}^{N_d} \log P(w_{nd}^* | \mathcal{A})$$

$$= \sum_{m=1}^M c_m^* \log P(w^* = m | \mathcal{A}) = \mathbf{c}^{*T} \log \hat{\beta}^*$$

$$(5) \quad P(\{w_{nd}^*\}_{n=1}^{N_d} | \mathcal{A}) = \frac{N_d!}{\prod_{m=1}^M c_m^*!} \prod_{n=1}^{N_d} P(w_{nd}^* | \hat{\beta}^*)$$

$$(6) \quad p(d) = \exp\left(-\frac{L(d)}{N_d}\right)$$

Table 1 contains the per-word perplexities for Document 2001 of the test set \mathcal{B} and the whole of \mathcal{B} , as well as the perplexity of \mathcal{B} when using a uniform multinomial. This is just equal to the number of distinct words in \mathcal{B} , as a perplexity of n implies a confidence the same as if an

```

1 log_posterior_frequencies = np.log(posterior_frequencies)
2 test_doc_log_probability = np.dot(test_doc_word_counts, log_posterior_frequencies)
3 test_doc_perplexity = np.exp(-1 * test_doc_log_probability / test_doc_total_words)

```

LISTING 3. Task C code excerpts

—	Document 2001	All Documents	Uniform Multinomial
Per-word Perplexity	4399.0	2697.1	6906

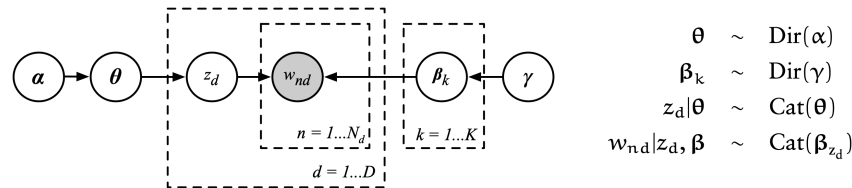
TABLE 1. Task C Perplexities

n -sided dice was choosing each word. Different documents have different perplexities due to the fact that they contain a different subset of the whole word set that is trained over. A document that has a higher proportion of words that were more common in the training set will have a higher log probability and therefore a lower (better) perplexity. Document 2001 has a much higher perplexity than \mathcal{B} as a whole, and therefore must contain a higher proportion of words that are scarcer in \mathcal{A} than the average.

4. TASK D: BAYESIAN MIXTURE OF MULTINOMIALS

Figure 3 shows an overview of the model used in Task D. $z_d \in 1, \dots, K$ is a latent variable denoting the topic of each document. The probability of a document being about topic k is $\theta_k \in \boldsymbol{\theta}$, which is initialised with a symmetric Dirichlet prior.

The python function `BMM` was provided which returns the test set perplexities and multinomial weights over words using gibbs sampling. The function was modified to return the number of documents assigned to each topic at each iteration by adding a new return array which contains a copy of `sk_docs` at each iteration, as seen in Listing 4.

FIGURE 3. Bayesian Mixture of Multinomials Model¹

```

1 sk_docs_history = np.zeros(
2     (num_mixture_components, num_iters_gibbs + 1), dtype=int
3 )
4 sk_docs_history[:, 0] = np.copy(sk_docs[:, 0])
5 # ...
6 for iteration in tqdm(range(num_iters_gibbs)):
7     # Do gibbs ...
8     sk_docs_history[:, iteration + 1] = np.copy(sk_docs[:, 0])

```

LISTING 4. Task D code excerpts

The model could then be trained using Gibbs Sampling over the training set using $\alpha = 10$ for the prior on mixture components and $\gamma = 0.1$ for the prior on words. The θ_k s can be calculated in the same way as β s earlier, shown in equation 7, although \mathbf{c} now is a vector of the topic counts.

$$(7) \quad \theta_k = \frac{\mathbf{c} + \alpha \cdot \mathbf{1}}{\sum_{l=1}^M c_l + \alpha}$$

Figure 4 shows the topic allocation progress for 3 different values of `np.random.seed`. For the most part, each reaches a stationary distribution, implying convergence of the model to a local minimum. However, different seeds yield different distributions, which are not explainable

¹CUED 4F13 Lecture notes (course site)

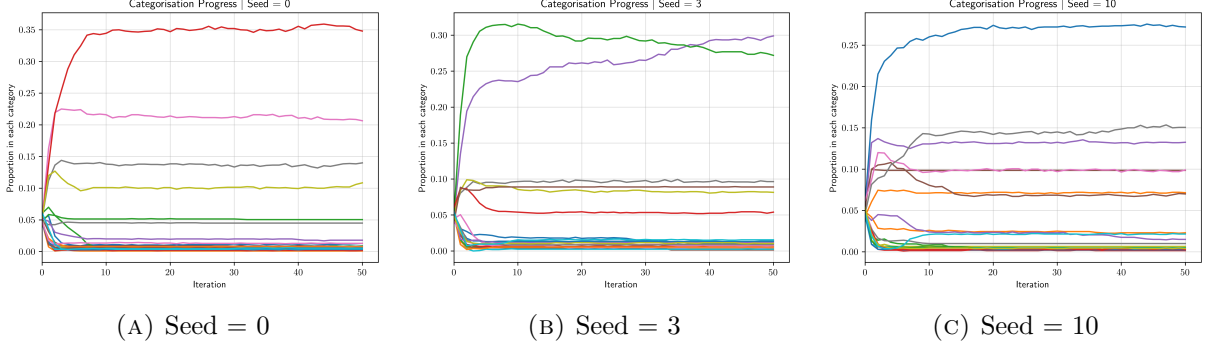
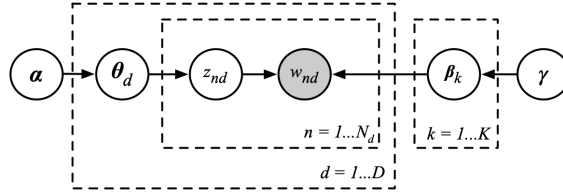


FIGURE 4. BMM Progress

by relabelling. Also, the vast majority of the 20 clusters are pretty much ignored (i.e. the model does not “believe” there are 20 distinguishable clusters). It cannot therefore be said that the algorithm is finding the same global optimum each time, and therefore is not truly exploring the “true” posterior. However, the test set perplexities for the three seeds are 2111, 2124, and 2116, which (averaged) is a $\sim 21\%$ improvement over the model in Tasks B and C. Therefore the model must work well to a certain extent.

5. TASK E: LATENT DIRICHLET ALLOCATION

In Latent Dirichlet Allocation (LDA), each word in a document can be drawn from a different topic, and therefore each document now has its own distribution over topics θ_d . Figure 5 shows the model, and the difference between LDA and BMM can be seen as there is now a θ_d for each document, rather than only one θ for the whole model. This increases the number of parameters significantly, making the need for Gibbs Sampling greater.

FIGURE 5. Latent Dirichlet Allocation Model²

The provided python function `LDA` was modified in a similar way to in Task D to return the topic posteriors at each Gibbs iteration, and also to return the word entropies for each topic at each iteration.

The model was run with parameters $\alpha = 0.1, \gamma = 0.1$ as before, giving a per-word perplexity of 1644, by far the lowest of all the models so far. This is as desired - we are increasing the complexity of the model but still decreasing the log likelihood of the test set - therefore the model is not overfitting to the training set.

Figure 6 shows the topic posteriors as a function of the Gibbs iteration for 3 different documents. The Gibbs sampler exhibits the *rich get richer* property - the more words in a document already assigned to a topic, the more a word is likely to be assigned to that topic. This explains why each document has one topic that dominates. It is possible that more than 50 iterations would be useful, as the topic allocations do not seem to have reached a stationary distribution. It was unfortunately not possible to run for more than 50 iterations to test this as the function was taking over 30 minutes to run.

The oscillations that persist may also be due to high similarity between multiple topics, which could be due to having too many topics, meaning a word has a similar probability of being drawn into either topic.

The word entropy can be calculated using equation 8. This is the expected uncertainty associated with each observation. Figure 7 shows the word entropies for each topic against

²See Footnote 1

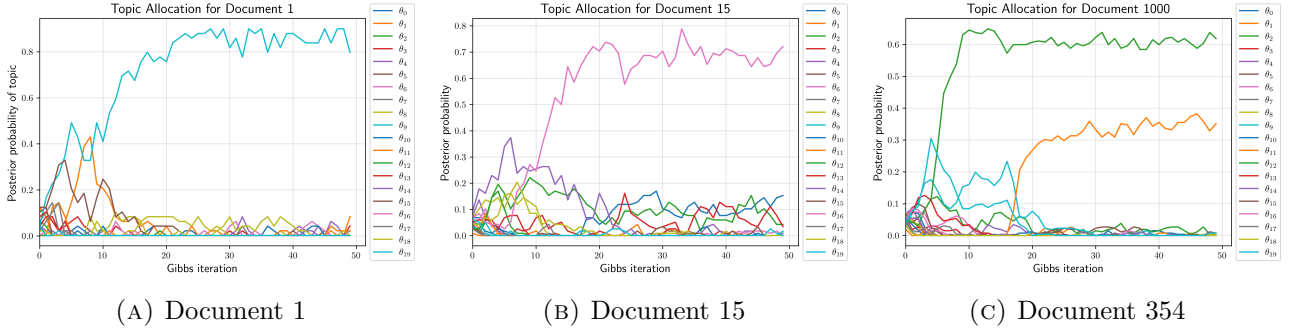


FIGURE 6. LDA Topic Progress

the Gibbs iteration. Initially, all the topics have the same entropy. After some iterations the entropy decreases, as the topics are reaching a distinct, non-uniform distribution - the topics are becoming more specific as training occurs, as expected. This means fewer words are required to “describe” the topic, or the topic has a smaller vocabulary - therefore the entropy decreases. The entropies reach fairly stable values after 50 iterations, implying that not many, if any, more iterations are required to train the model any better.

$$(8) \quad H_k = - \sum_{m=1}^M \beta_{km} \log(\beta_{km})$$

The natural logarithm is used instead of the standard base 2 logarithm as that allows for easier comparison with the perplexity. Using the natural logarithm results in units of nats, not bits. The log-perplexity is a comparable unit to the entropy, as they will be the same if measuring an infinitely long document on a single topic. The log-perplexity of the test set is $\log 1644 = 7.40$, which is higher than all the topic entropies, as expected as the test set is a blend of all the topics, and therefore less specific than any single topic.

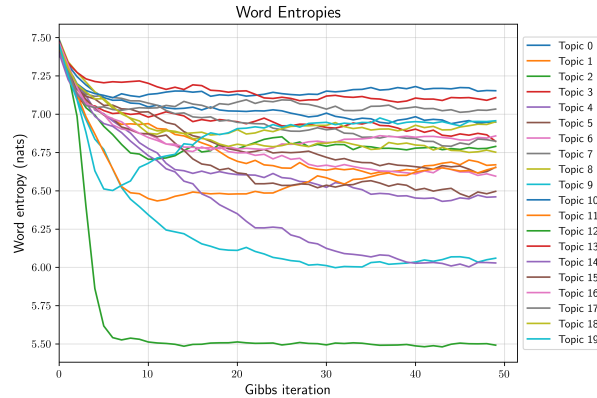


FIGURE 7. Entropy Evolution