
Customer Feedback Topic Modelling Using Online Latent Dirichlet Allocation

Yuhou Zhou

Department of Computer Science & Electrical Engineering
Jacobs University Bremen
Bremen, Germany
Corporate Information Management
CEWE Stiftung & Co. KGaA
Oldenburg, Germany
`yu.zhou@jacobs-university.de`
`yuhou.zhou@cewe.de`

Abstract

I develop a topic modelling system¹ for extracting topics from customer feedback. The system includes three parts which are preprocessing, modelling, and visualization. Apache PySpark helps to read large volume data, and Python spaCy is used to run preprocessing pipeline. Online Latent Dirichlet Allocation is the selected topic model, and the system uses Gensim implementation. Topic coherence serves as the evaluation metric for choosing the optimal topic number. LDavis presents topic modelling results interactively and shows the topic distribution over two principal components. This project can serve as a base for further generalizing to extract topics from different data sources.

1 Introduction

Customer feedback is information provided by clients about their general experience with a product or a service. At CEWE, one method to collect such feedback is asking customers to complete online survey.

By collecting customer feedback, the extracted information can act as a guide to improve products. Feedback gives the company the clue of customers' general experience and the satisfaction rate. Except the description provided by sellers, the publicly posted feedback is useful information to other consumers, imparting them an overview how the product really is. This information also helps managers to make correct business decisions.

The current way of processing user feedback requires that working staff reads the feedback and later extracts topics. This process can be costly and inefficient, because of manual intervention. At one time, people can only process a small batch of text from a recent period. Thus, the information, which concerning long lasting but less frequently appeared issues, may lost.

Latent Dirichlet Allocation (LDA) can automate this topic extraction process. LDA is a generative probabilistic model for discrete data, and it is especially popular in the field of topic modelling. In this project, I automate the topic extraction process utilizing Online Latent Dirichlet Allocation with the optimal number of topics found by topic coherence.

¹Available on gitlab.com/yuhouzhou/cftm

2 Related Work

Discovering “topics” from a collection of documents has been a long-standing problem in natural language understanding. Latent Semantic Analysis (LSA) acts as the base of recent topic models [1]. Hofmann proposed Probabilistic Latent Semantic Indexing (pLSA) [2], bringing topic models from deterministic to probabilistic. Blei et al. (2003) described Latent Dirichlet Allocation (LDA) which is a Bayesian alternative to pLSA [3]. LDA makes use of word co-occurrence information in documents, to generate document-topic distribution. Hoffman et al. (2010) developed an online variational Bayes algorithm for LDA [4]. This online version of LDA was demonstrated that it has advantages over batch LDA in terms of model accuracy and the modelling time. However, the lack of word co-occurrence in short texts prevents LDA yielding satisfactory results [5]. Such short texts are news headlines, tweets, user feedback, questions and answers, and so on. One strategy to overcome this data sparsity problem is to concatenate a set of short texts according to available metadata [6, 7, 8]. For example, aggregating texts by timestamps, locations, or users.

Perplexity was used to evaluate the accuracy of topic models [9]. However, Chang et al (2009) found that perplexity does not always correlated with the interpretability of topics [10]. Topic coherence C_{UCI} was first proposed by Newman et al.(2010)[11], by using pairwise pointwise mutual information (PMI) between the topic words. C_{UCI} was proven highly correlated with human judgment. Röder et al. (2015) proposed topic coherence C_V which was identified that it outperformed other topic coherence [12].

3 Online Latent Dirichlet Allocation

Latent Dirichlet Allocation is an unsupervised machine learning model [3]. It identifies latent topics in corpus. Each document in the corpus is represented as probability distribution over some other topics. All the topics are represented as a probability distribution over some words.

Theoretically, a document w in the corpus D is associated with a multinomial distribution over T topics; T is the predefined number of topics, and document is denoted by $w = (w_1, w_2, \dots, w_N)$ where w_n is the n th word in the sequence. The multinomial distribution is denoted as θ . One topic is similarly associated with a multinomial distribution ϕ over words w_n .

θ and ϕ whose hyper-parameters are α and β respectively have Dirichlet prior. For every word in one document w , a topic z is sampled consequently. This generative process is repeated N times, where N is the number of words in the document [7].

A k-dimensional Dirichlet Random variable has the following probability density:

$$p(\theta|\alpha) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \dots \theta_k^{\alpha_k-1} \quad (1)$$

where $\Gamma(x)$ is the Gamma function. The probability of a corpus:

$$p(D|\alpha, \beta) = \prod_{d=1}^M \int p(\theta_d|\alpha) \left(\prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn}|\theta) p(w_{dn}|z_{dn}, \beta) \right) d\theta_d \quad (2)$$

The probabilistic graphical model representation of LDA is in Figure 1. As shown in the Figure 1, LDA consists of three layers. The parameters α and β are on the level of corpus. θ_d are variables on the level of document. z_{dn} and w_{dn} are variables on the level of word.

Batch LDA requires to feed the entire corpus into each iteration when the model is being trained, so in some circumstance LDA does not perform well, such as when the datasets are large, or when new data comes in stream to the training process [4].

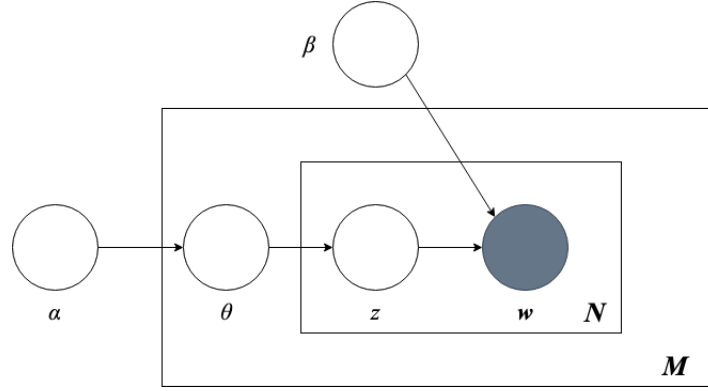


Figure 1: Graphical model representation of LDA

Hoffman et al. (2010) proposed an online variational Bayes (VB) for LDA, instead of using "batch" VB or Markov Chain Monte Carlo (MCMC) as in batch LDA. VB has the advantage of memory efficiency comparing to MCMC, and it also empirically faster than MCMC. In his experiments, online LDA also converges much faster than the batch algorithm [4].

4 Topic Coherence

After perplexity, a more accurate metric to evaluate the coherence of topic modelling was proposed by Newman et al [11]. This measure takes N top words from a topic and sum a confirmation measure over all words. This topic coherence is based on pointwise mutual information (PMI) and calculated as:

$$C_{UCI} = \frac{2}{N \cdot (N - 1)} \sum_{i=1}^{N-1} \sum_{j=i+1}^N PMI(w_i, w_j) \quad (3)$$

$$PMI(w_i, w_j) = \log \frac{P(w_i, w_j) + \epsilon}{P(w_i) \cdot P(w_j)} \quad (4)$$

Mimno et al. use a different confirmation measure which is asymmetric to replace PMI. It is calculated as:

$$C_{UMass} = \frac{2}{N \cdot (N - 1)} \sum_{i=2}^N \sum_{j=1}^{i-1} \log \frac{P(w_i, w_j) + \epsilon}{P(w_j)} \quad (5)$$

Röder et al (2015) proposed a framework which can generate millions of topic coherence metrics, by combining elementary components [12]. At present, they already evaluated 200,000 coherences. Four components consist of this framework. In the first component, the word set t is divided into a set of pairs of word subsets S , and this stage is denoted as \mathcal{S} . In the second stage, word probabilities P are computed by using a given reference corpus. This stage is denoted by \mathcal{M} . The third stage \mathcal{P} uses S and P to calculate the agreement ϕ of the word pairs in S . In the final stage Σ , those values are aggregated to a single coherence value c . The configuration space is the cross product of the four sets:

$$\mathcal{C} = \mathcal{S} \times \mathcal{M} \times \mathcal{P} \times \Sigma \quad (6)$$

Röder et al evaluated that the topic coherence C_V outperforms other topic coherence [12]. C_V combines the indirect cosine measure with the NPMI, which is normalized PMI, and the boolean sliding window, which determines word counts by a sliding window.

5 LDAvis

Sievert et al. (2014) presented a web-based visualization for Latent Dirichlet Allocation [13]. It projects high dimensional topic space into two principal components. To interpret a topic, one typically examines a ranked set of term by its *probability*.

Instead presenting words pure by their *probability* in a topic, LDAvis presents words according their *relevance* to the topic, which can highlight the words with high *exclusivity* or *lift* [14] and give a possibility to reduce the rank of universal words. *relevance* is defined as:

$$r(w, k|\lambda) = \lambda \log(\phi_{kw}) + (1 - \lambda) \log \frac{\phi_{kw}}{p_w} \quad (7)$$

where λ determines the weight given to the *probability* of word w under topic k relative to its *lift*. The higher λ gives more weight to its *probability* and less to its *lift*. ϕ is estimated in Online LDA using Online Variational Bayes methods, and p_w is from the empirical distribution of the corpus.

6 Data Source

The original data files used for this project are two parquet files of the overall size 1.47 GB. The parquet file in customer_feedbacks/ contains 101 columns, and the one in customer_feedbacks_cat/ contains 7 columns. The content of customer feedback is stored in the column ANTWORT_WERT, There are 1,917,490 feedbacks including duplicates since 14/04/2012. The sentiment of the feedback is stored in the column STIMMUNG.

7 System Architecture and Usage

Under the project root directory, there are cftm/ directory, environment.yml file and README.md file.

environment.yml files specifies the environment dependencies of the system. We can use conda to recreate the system environment. README.md provide brief introduction of the system.

The source code of the system is under cftm/cftm/ directory. It contains three python files which are pipeline.py, cftm_parser.py, preprocessing.py and path.yml file.

pipeline.py is the entry of the system, passing arguments to this file in order to control the behavior of the system. cftm_parser.py and preprocessing.py contains modules used by pipeline. cftm_parser.py extracts specified data from data source file; preprocessing.py executes standard nlp preprocessing operations, eventually converting texts to the object which Gensim requires. path.yml file contains all the paths to store system output.

7.1 cftm_parser.py

cftm_parser.py contains one method `parquet_transform()` to transform data in parquet files to the data format for preprocessing and modelling.

`parquet_transform(path1, path2, n=-1)`

parameters:

- path1 (str) - the path of the first parquet file
- path2 (str) - the path of the second parquet file
- n (int, optional) - the number of observations to be transformed. When the passed integer is smaller than 0, all observations from data source will be transformed. The default value is -1.

returns:

- `pandas.dataframe` - include the data frame with desired columns and the specified number of observations

cftm_parser makes use of Apache Spark, and Spark SQL to import data from parquet file. Since the size of the parquet files are in the level of gigabytes, directly using method `pandas.read_parquet()` to extract pandas data frame requires large volume of memory. This is sometimes unrealistic when only a personal computer is available for running the system.

An alternative way to import the data is using `spark.read.parquet()`. Instead of fitting all the data to memory, Spark's operators spill data to disk, which make it possible to run over any size of data. `spark.read.parquet()` reads a parquet file from a specified path and return a `spark.dataframe` object. Writing queries by Spark SQL makes the data manipulation the same as the process over relational database. After the query the size of the data drops significantly. Sequentially, duplicated observations are dropped and `spark.dataframe` is transformed to `pandas.dataframe` for preprocessing procedures.

7.2 preprocessing.py

`preprocessing.py` utilizes python spaCy to run the standard text preprocessing pipeline, which includes tokenization, lemmatization, lowercasing, stripping white spaces, and removing stopwords and punctuations. A token is the name for a sequence of characters that we want to treat as a group, i.e. word, punctuation symbol, white space, etc. Tokenization means to convert a text to a collection of tokens. Lemmatization assigns the base forms of words. For example, assign "run" to "running", "ran", or "be" to "is", "am", "were", etc. Stopwords are the words which appear with high-frequency, such as "also", "the", "to", etc. They too ubiquitously exist in texts, and therefore contain little information about the text. Stopwords and punctuations are removed to only keep words which are distinct.

spaCy² is a Python/Cython natural language processing package which currently offers fastest syntactic parser in the word. Choi et al. (2015) conducted an evaluation research proved the outstanding performance of spaCy [15]. The high speed of spaCy makes it especially suitable to processing large volume of text data. spaCy currently supports 8 languages, German is one of them. This makes spaCy be able to be used in this project.

Four methods are in this script. `text_preproc_maker()`, `text_aggregation()`, `gensim_prep()`, and `preprocessor()`.

`text_preproc_maker(stopwords, language='de')`

A currying function which returns a new function `text_preproc(sentence)` to run the text preprocessing pipeline, including tokenization, lemmatization, lowercasing, stripping white spaces, and removing stopwords and punctuations, for the specified language.

parameters:

- stopwords (list of str) - a list of stopwords
- language (str, optional) - the possible inputs are the languages which spaCy supports, such as 'de', 'en', etc. Before use the specified language, first install the spaCy required language, using command `pip install spacy && python -m spacy download en`, here 'en' can be replaced with the language you would like to install.

returns:

- `text_preproc object` - a function takes a str as the parameter, containing the pipeline of preprocessing.

`text_aggregator(df_pd, metadata='DATE', min_len=-1)`

parameters:

- df_pd (pandas.dataframe) - a Pandas data frame containing text data and metadata.
- metadata (str, optional) - the column name of the metadata which is used to aggregate texts. Legal options are 'DATE' and 'SENTIMENT'.

²Available on github.com/explosion/spaCy

- `min_len` (int, optional) - when the metadata is 'DATE', this argument decides the minimal length of the aggregated text. The default is -1, which means no aggregation but using the original texts.

returns:

- *list of lists of str* - the text data after aggregation.

gensim_prep(texts)

This method converts list of lists string to gensim dictionary and corpus objects.

parameters:

- `texts` (list of lists of str) - text data.

returns:

- *list of lists of str* - the text data.
- *gensim.corpora.dictionary.Dictionary* - A mapping between words and their integer ids.
- *list of lists of (int, int)* - Bag of words representation of the documents.

preprocessor(df_pd, stopwords, language='de', text='TEXT', metadata=None, min_len=-1)

This method wraps up preproc, text aggregator, and gensim_prep to compose a full preprocessing pipeline. It returns text data, gensim dictionary, and gensim bag of words corpus. See the information of the arguments above.

7.3 pipeline.py

pipeline.py is the entry of the system. It comprises three parts – preprocessing, modelling, and visualization. Preprocessing part uses the function *preprocessor()* in preprocessing.py. For modelling, gensim implemented Online LDA is used. pyLDavis module takes the training results from Online LDA, and produce a html file which supports interactive visualization.

7.3.1 Arguments

This program takes 9 parameters. *-pipeline* is a positional argument, and others are optional arguments.

-pipeline is a positional parameter, requiring 3 integer numbers, which switches on or off preprocessing, modelling, and visualization respectively. If one of the numbers is 0, the corresponding function is turned off. If it is none 0, the corresponding function is turned on. For example, if the parameter is *1 1 1*, then the pipeline will be preprocessing, modelling, and visualization. If the parameter is *0 1 1*, then the pipeline will only have the procedures of modelling and visualization, when the training data is read from a pickle file, which is either got from an archive or is copied to the required path. The parameter *0 0 1* has the same logic. *1 0 0* and *1 1 0* also are the legal combinations. Designing the system in this way makes the pipeline flexible and fault-tolerant. Since the data volume is big, user can choose which part of the pipeline he would like to run to avoid repeated work. Every part of the pipeline produces a pickle file to archive the outcome, so even if the program is stopped during a run, the results are saved and can be continued next run.

-path_file or *-p* is an optional parameter, requiring a str value. The default value is *./path.yml*. It accepts a path of the yml file in which all output paths are indicated. The yml file uses a dictionary to store paths, and the keys are *parquet_path1*, *parquet_path2*, *archive_fore_path*, *data_path*, *model_path*, *pic_path*, *html_path*, *log_path*.

-observation_n or *-o* indicates the number of observations used in the program, requiring an integer. The default value is -1, which means use all observations.

-agg_metadata or *-m* indicates the metadata which is used during the aggregation session. The default value is 'DATE', and the other legal value is 'SENTIMENT'.

-agg_length or *-l* indicates the minimal length of every document after an aggregation session, if the metadata 'DATE' is used. It requires an integer number, and the default value is -1, which means the aggregation is aborted.

`-topic_range` or `-r` indicates the topic range during the modelling session. It requires 3 integer numbers, meaning the minimal number of topics, the maximal number of topics, and the step. The default value is `1 10 1`, so the modelling session will generate 10 LDA models with topic numbers from 1 to 10 topics.

`-chunksize` or `-c` is a parameter to Online LDA model. Online LDA algorithm update the model every chunksize. For example, if a corpus has 10000 documents, and the chunksize is set to 1000, then Online LDA will update the model 10 times through one complete pass of the corpus. The default value is 2000.

`-iteration` or `-i` indicates the number of iterations during optimization, requiring an integer number. The number of iterations should be high enough to make the model converge. However, fewer iterations make the modelling time shorter. The default value is 50.

`-seed` or `-s` indicates the seed passed to LDA model, make the whole process replicable. The default value is 1.

An example command can be `python pipeline.py 1 1 1 -p ./path.yml -o 40000 -m DATE -l 15 -r 10 50 5 -c 800 -i 400 -s 2`, which means: run the complete pipeline (preprocessing, modelling, and visualization). In preprocessing, 40000 observations are included, and the system uses 'DATE' column as metadata to aggregate documents. After aggregation, every document has minimum length of 15 tokens. In modelling session, LDA models with topic numbers 10, 15, 20, ..., 50 are generated. 800 observations as a chunk is used for training, in this case, one full pass of corpus will update models $40000/800 = 50$ times. There are maximum 400 iterations through the corpus when inferring the topic distribution of the corpus. All the output files, such as processed training data, log files, lda models, topic coherence models, visualizations, etc, are saved in the paths specified in `./path.yml`. The seed is set as 2.

For every run of the program, parameters are saved in `commandline_args.yml` file in the output folder and archive folders.

7.3.2 Modelling

gensim³ is an open-source library for unsupervised topic modelling and natural language processing. This software framework was proposed by Řehůřek in 2010, featuring with its clarity, efficiency and scalability [16]. gensim supports parallel computing and is scalable through a cluster, which accelerates training process significantly. For this customer feedback topic extraction system, gensim Online LDA model is chosen.

If in one program run, the preprocessing stage is turned on. The processed data will be used in the model training session. If not, the training data will be read from previous pickled training data.

The modelling session generates LDA models with the number of topics passed to `-topic_range` or `-r` argument. For every model, its topic coherence is also generated. In this project, topic coherence C_V is used, because it often outperforms other topic coherence [12]. After one model is trained and its topic coherence is evaluated, it will be dumped as `.pickle` file. This is a full iteration of modelling, and next iteration the LDA model is passed with a topic number which is `step-size` higher. In the end of every modelling iteration, the model and its topic coherence will be added to previous `.pickle` file.

7.3.3 Visualization

Visualization includes two parts. The first part is generating scatterplot, topic coherence of the model C_V against topic numbers. The higher the score of topic coherence often means better interpretability of the topics. The topic number with the highest topic coherence is marked red on the scatterplot. The second part is visualize the topics of the model with the highest topic coherence score, using python package pyLDAvis. The html file is dumped.

³Available on github.com/RaRe-Technologies/gensim

8 Experiments

The system is designed to be flexible and easy to use. Users can include an observation number depending on the time they are able to afford and switch on or off parts of the pipeline as they need. To archive the best modelling result, users should include observations as many as possible, i.e. set -1 to the argument $-o$. In some situations, users may not want to include all the observations from data set. Such as to reduce the training time, to test the system, or to know past few years topics, etc, so in this case users can specify the number of observations to be included.

When doing topic modelling on short tests, such as user feedback, questions and answers, etc, LDA models usually have unsatisfactory results in which the interpretability of topics suffers from lack of word co-occurrence. In order to solve short texts topic modelling problem, different aggregation methods can be applied to the original texts, composing a long document. Aggregation strategies usually make use of available metadata, such as text date, sentiment, user source, etc [6, 7, 8].

To generate an optimal model, users only need to set $-chunksize$ sufficiently small to guarantee high update times; set $-iteration$ big enough to make models converge; set $topic_range$ cover a wide interval of topic numbers. Experience of setting these arguments can be gained by inspecting data sets and repeated experiments.

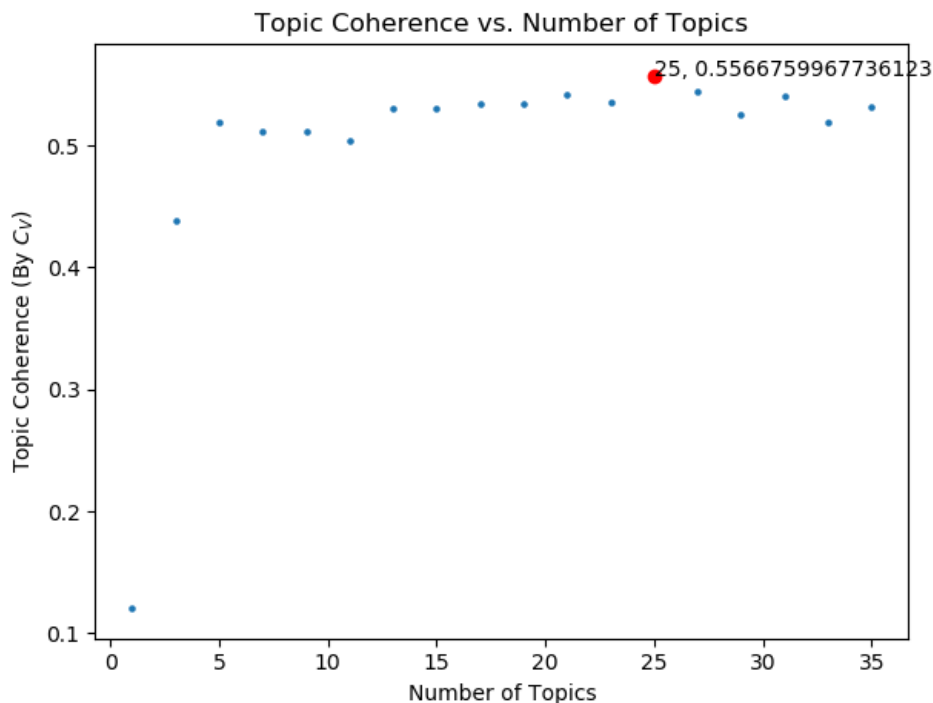


Figure 2: Topic Coherence vs. Number of Topics

We yield the best LDA model with 25 topics using 1.9 million customer feedback, by using following parameters:

`python pipeline.py 1 1 1 -c 8000 -i 600 -r 1 36 2` (8)

The scatterplot of topic coherence against number of topics is shown in Figure 2, and the screenshot of the topic visualization is in Figure 3. The top 5 topics are shown below:

Topic 1: *filiale, bestellung, erhalten, abholen, lieferung, liefern, dm, kalender, weihnachten, ...*
 Topic 2: *bild, rand, abschneiden, leinwand, weiß, zuschneiden, automatische, mm, kopf, dunkel, ...*
 Topic 3: *möglichkeit, hintergrund, seite, layout, cliparts, einfügen, auswahl, gestalten, software, ...*

Topic 4: *fotobuch, cewe, empfehlen, begeistern, kunde, weiterempfehlen, freund, zufrieden, ...*
 Topic 5: *schnellen, lieferung, bearbeitung, zuverlässig, preiswert, zügig, prompt, einfach, unkomplizierte, unproblematisch, ...*

Interpreting the word context, we can know the ideas of these topics. Topic 1 is about the product delivery and pickup process. Topic 2 is about the cropping and the color quality of printed photos. Topics 3 is about the the photo book and calendar design in CEWE software. Topic 4 is about the satisfaction and recommendation to the products. Topic 5 stresses the quick speed of delivery and the easy delivery process.

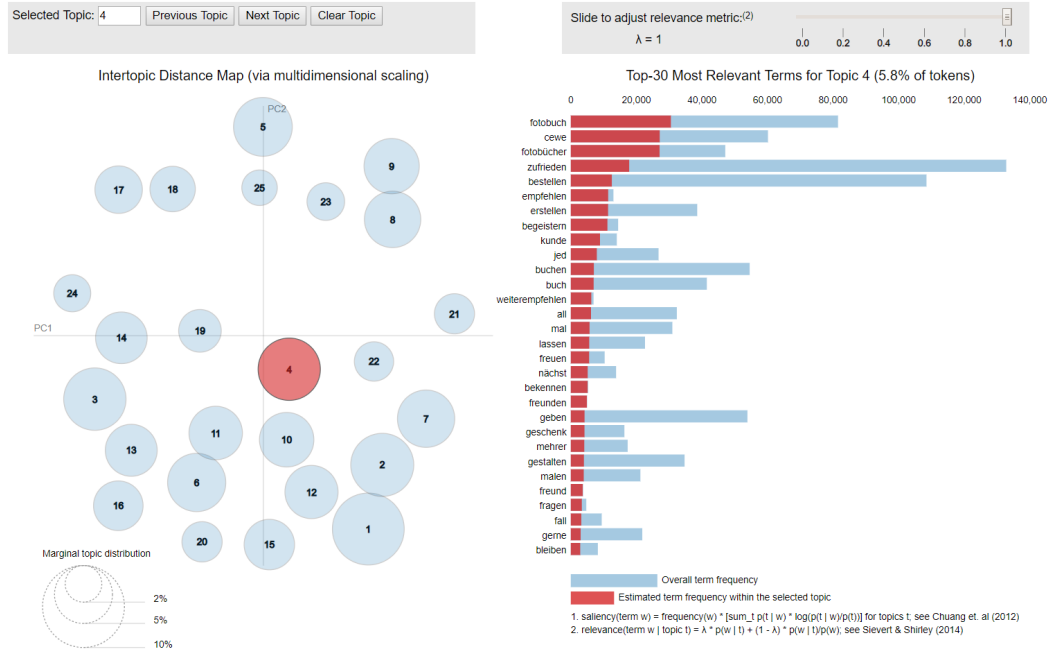


Figure 3: Screenshot of the Modelling Visualization

9 Discussion

Since the topic modelling is an interesting task in many situations, such as extracing topics from articles, get top topics from tweets of different users, etc. It will be very helpful, if the system can be further generalized by providing an interface where the system gets different parsers. The ability to directly connect to databases and run queries on them should also be included in the system in the future.

The evaluation metrics can be improved. Some results from different topic coherence indexes conflict with each other. For example, topic coherence generated from C_{UMass} [17] and C_V [12] are often negatively correlated. Topic coherence thus is not so reliable to be the single metric to find the best topic numbers. To improve the consistence of topic coherence indexes, further studies should be conducted.

Currently, the system is using terminal as the interface with users. The system can be more approachable to non-technical users by providing a better user interface, utilizing tools such as Semantic UI to build a graphical user interface.

The system is purely written in Python, and uses conda to manage its environment. However, the process of using conda to rebuild the exactly same environment is not that easy. Some dependency packages in the environment are operating system dependent, for example *libgfortran* is only available on macOS and linux systems. If the user is using a Windows system, the package cannot be found. This can make the system deployment cumbersome. Containerizing the system using Docker can

be a possible solution to this problem. Containerized systems can be easily distributed to different operating systems, and the dependencies can be preinstalled.

Aggregation can improve the interpretability of topics. It was proved in our experiments. However, after applying aggregation strategy in our project, granularity of topics reduced. This is maybe caused by the dramatic decrease of training observations after aggregation, because the original average token number of all customer feedback is around 9. The mechanism still need to be studied and the aggregation strategy can be further refined to improve the performance of the system.

Gensim LDA provides distributed computing ability. Distributed Online LDA makes use of a computer cluster and distribute tasks to different nodes while modelling. This feature can significantly reduce modelling time and especially suitable for the situation where the data volume are too big to process on single computer, or it takes too long to finish modelling. The configuration to a computer cluster should be made in the future.

10 Summary

In this paper, I implement a customer feedback topic modelling system using Online Latent Dirichlet Allocation (LDA). The system includes the full data pipeline from text preprocessing to modelling and final to topic visualization. The difficulties of picking the topic number is solved by evaluating topic coherence scores of models. The final model is presented interactively. A full pickle and archive mechanism is developed, making the system fault-tolerant, and past results can be reviewed readily. Future work could involve generalizing the system to a broader use, improving evaluation metrics, developing a better graphical user interface, containerizing the system, and scaling the system.

Acknowledgments

This work is supervised by Sadok Ben Yahya and Fatih-Mehmet Inel at CEWE Stiftung & Co. KGaA, and mentored by professor Adalbert F.X. Wilhelm at Jacobs University Bremen. This project is established on the Customer Feedback Classification System of Corporate Information Management, CEWE Stiftung & Co. KGaA.

References

- [1] S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman, "Indexing by latent semantic analysis," *Journal of the American Society for Information Science*, vol. 41, no. 6, pp. 391–407, Sep. 1990. [Online]. Available: <http://doi.wiley.com/10.1002/%28SICI%291097-4571%28199009%2941%3A6%3C391%3A%3AAID-ASII%3E3.0.CO%3B2-9>
- [2] T. Hofmann, "Probabilistic Latent Semantic Indexing," *ACM SIGIR Forum*, vol. 51, no. 2, p. 8, 2017.
- [3] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet Allocation," *Journal of Machine Learning Research*, vol. 3, no. Jan, pp. 993–1022, 2003. [Online]. Available: <http://www.jmlr.org/papers/v3/blei03a.html>
- [4] M. Hoffman, F. R. Bach, and D. M. Blei, "Online Learning for Latent Dirichlet Allocation," in *Advances in Neural Information Processing Systems 23*, J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, Eds. Curran Associates, Inc., 2010, pp. 856–864. [Online]. Available: <http://papers.nips.cc/paper/3902-online-learning-for-latent-dirichlet-allocation.pdf>
- [5] C. Li, H. Wang, Z. Zhang, A. Sun, and Z. Ma, "Topic Modeling for Short Texts with Auxiliary Word Embeddings," in *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval - SIGIR '16*. Pisa, Italy: ACM Press, 2016, pp. 165–174. [Online]. Available: <http://dl.acm.org/citation.cfm?doid=2911451.2911499>
- [6] J. Weng, E.-P. Lim, J. Jiang, and Q. He, "TwitterRank: finding topic-sensitive influential twitterers," in *Proceedings of the third ACM international conference on Web search and data*

- mining - WSDM '10. New York, New York, USA: ACM Press, 2010, p. 261. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1718487.1718520>
- [7] L. Hong and B. D. Davison, "Empirical study of topic modeling in Twitter," in *Proceedings of the First Workshop on Social Media Analytics - SOMA '10*. Washington D.C., District of Columbia: ACM Press, 2010, pp. 80–88. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1964858.1964870>
 - [8] R. Mehrotra, S. Sanner, W. Buntine, and L. Xie, "Improving LDA topic models for microblogs via tweet pooling and automatic labeling," in *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval - SIGIR '13*. Dublin, Ireland: ACM Press, 2013, p. 889. [Online]. Available: <http://dl.acm.org/citation.cfm?doid=2484028.2484166>
 - [9] H. M. Wallach, I. Murray, R. Salakhutdinov, and D. Mimno, "Evaluation methods for topic models," in *Proceedings of the 26th Annual International Conference on Machine Learning - ICML '09*. Montreal, Quebec, Canada: ACM Press, 2009, pp. 1–8. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1553374.1553515>
 - [10] J. Chang, S. Gerrish, C. Wang, J. L. Boyd-graber, and D. M. Blei, "Reading Tea Leaves: How Humans Interpret Topic Models," in *Advances in Neural Information Processing Systems 22*, Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta, Eds. Curran Associates, Inc., 2009, pp. 288–296. [Online]. Available: <http://papers.nips.cc/paper/3700-reading-tea-leaves-how-humans-interpret-topic-models.pdf>
 - [11] D. Newman, J. H. Lau, K. Grieser, and T. Baldwin, "Automatic evaluation of topic coherence," p. 9.
 - [12] M. Röder, A. Both, and A. Hinneburg, "Exploring the Space of Topic Coherence Measures," in *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining - WSDM '15*. Shanghai, China: ACM Press, 2015, pp. 399–408. [Online]. Available: <http://dl.acm.org/citation.cfm?doid=2684822.2685324>
 - [13] C. Sievert and K. Shirley, "LDAvis: A method for visualizing and interpreting topics," in *Proceedings of the Workshop on Interactive Language Learning, Visualization, and Interfaces*. Baltimore, Maryland, USA: Association for Computational Linguistics, 2014, pp. 63–70. [Online]. Available: <http://aclweb.org/anthology/W14-3110>
 - [14] J. M. Bischof and E. M. Airoidi, "Summarizing topical content with word frequency and exclusivity," p. 8.
 - [15] J. D. Choi, J. Tetreault, and A. Stent, "It Depends: Dependency Parser Comparison Using A Web-based Evaluation Tool," in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Beijing, China: Association for Computational Linguistics, Jul. 2015, pp. 387–396. [Online]. Available: <http://www.aclweb.org/anthology/P15-1038>
 - [16] R. Řehůřek and P. Sojka, "Software Framework for Topic Modelling with Large Corpora," in *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. Valletta, Malta: ELRA, May 2010, pp. 45–50, <http://is.muni.cz/publication/884893/en>.
 - [17] D. Mimno, H. M. Wallach, E. Talley, M. Leenders, and A. McCallum, "Optimizing semantic coherence in topic models," p. 11.