



Text Summarization

Text summarization refers to the task of compressing a relatively large amount of text data or a long text article into a more concise form for easy digestion. It is obviously very important for text data access, where it can help users see the main content or points in the text data without having to read all the text. Summarization of search engine results is a good example of such an application. However, summarization can also be useful for text data analysis as it can help reduce the amount of text to be processed, thus improving the efficiency of any analysis algorithm.

However, summarization is a non-trivial task. Given a large document, how can we convey the important points in only a few sentences? And what do we mean by “document” and “important”? Although it is easy for a human to recognize a good summary, it is not as straightforward to define the process. In short, for any text summarization application, we’d like a semantic compression of text; that is, we would like to convey essentially the same amount of information in less space. The output should be fluent, readable language. In general, we need a purpose for summarization although it is often hard to define one. Once we know a purpose, we can start to formulate how to approach the task, and the problem itself becomes a little easier to evaluate.

In one concrete example, consider a news summary. If our input is a collection of news articles from one day, a potentially valid output is a list of headlines. Of course, this wouldn’t be the entire list of headlines, but only those headlines that would interest a user. For a different angle, consider a news summarization task where the input is one text news article and the output should be one paragraph explaining what the article talks about in a readable format. Each task will require a different solution.

Summarizing retrieval results is also of particular interest. On a search engine result page, how can we help the user click on a relevant link? A common strategy is to highlight words matching the query in a short snippet. An alternative approach would be to take a few sentences to summarize each result and display the short

summaries on the results page. Using summaries in this way could give the user a better idea of what information the document contains before he or she decides to read it.

Opinion summarization is useful for both businesses and shoppers. Summarizing all reviews of a product lets the business know whether the buyers are satisfied (and why). The review summaries also let the shoppers make comparisons between different products when searching online. Reviews can be further broken down into summaries of positive reviews and summaries of negative reviews. An even more granular approach described in [Wang et al. \[2010\]](#) and [Wang et al. \[2011\]](#) and further discussed in Chapter 18 uses topic models to summarize product reviews relating to different aspects. For hotel reviews, this could correspond to service, location, price, and value. Although the output in these two works is not a human-readable summary, we could imagine a system that is able to summarize all the hotel reviews in English (or any other language) for the user.

In this chapter, we overview two main paradigms of summarization techniques and investigate their different applications.

16.1 Overview of Text Summarization Techniques

There are two main methods in text summarization. The first is **selection-based** or **extractive summarization**. With this method, a summary consists of a sequence of sentences selected from the original documents. No new sentences are written, hence the summary is *extracted*. The second method is **generation-based** or **abstractive summarization**. Here, a summary may contain new sentences not in any of the original documents. One method that we explore here is using a language model. Previously in this book, we've used language models to calculate the likelihood of some text; in this chapter, we will show how to use a language model in reverse to generate sentences. We also briefly touch on the field of **natural language generation** in our discussion of abstractive techniques.

Following the pattern of previous chapters, we then move on to evaluation of text summarization. The two methods each have evaluation metrics that are particularly focused towards their respective implementation, but it is possible to use (e.g.) an abstractive evaluation metric on a summary generated by an extractive algorithm. Finally, we look into some applications of text summarization and see how they are implemented in real-world systems.

Text summarization is a broad field and we only touch on the core concepts in this chapter. For further reading, we recommend that the reader start with [Das and Martins \[2007\]](#), which provides a systematic overview of the field and contains much of the content from this chapter in an expanded form.

16.2 Extractive Text Summarization

Information retrieval-based techniques use the notion of sentence vectors and similarity functions in order to create a summarization text. A sentence vector is equivalent in structure to a document vector, albeit based on a smaller number of words. Below, we will outline a basic information retrieval-based summarization system.

1. Split the document to be summarized into sections or passages.
2. For each passage, “compress” its sentences into a smaller number of relevant (yet not redundant) sentences.

This strategy retains coherency since the sentences in the summary are mostly in the same order as they were in the original document.

Step one is portrayed in Figure 16.1. The sentences in the document are traversed in order and a normalized, symmetric similarity measure (see Chapter 14) is applied on adjacent pairs of sentences. The plot on the right-hand side of the figure shows the change in similarity between the sentences. We can inspect these changes to segment the document into passages when the similarity is low, i.e., a shift in topic occurs. An alternative approach to this segmentation is to simply use paragraphs if the document being operated on contains that information, although most of the time this is not the case. This rudimentary partitioning strategy is a task in

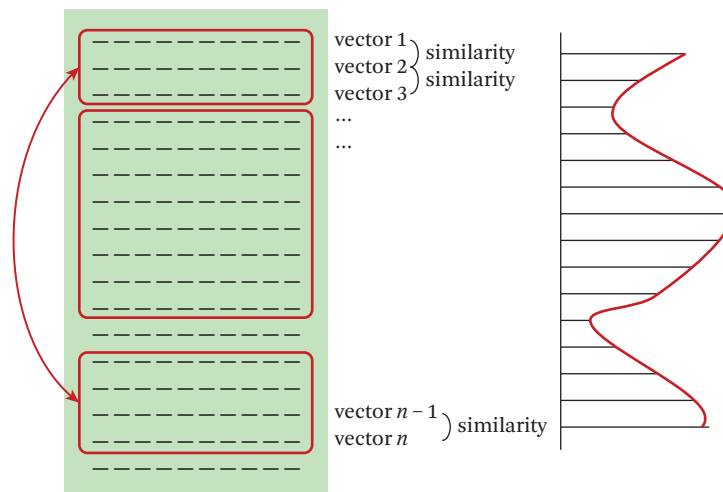


Figure 16.1 Segmenting a document into passages with a similarity-based discourse analysis.

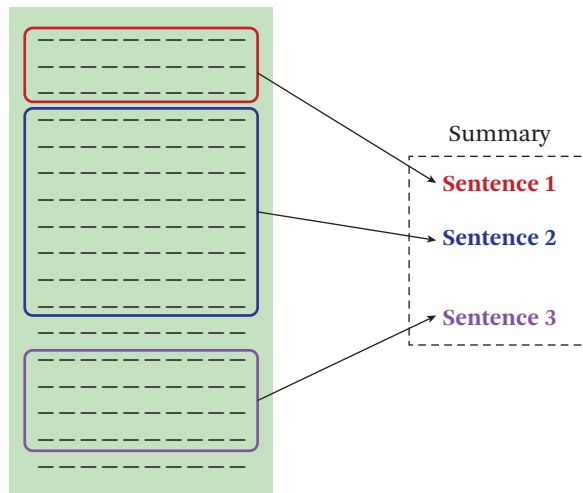


Figure 16.2 Text summarization using maximum marginal relevance to select one sentence from each passage as a summary.

discourse analysis (a subfield of NLP). Discourse analysis deals with sequences of sentences as opposed to only one sentence.

Now that we have our passages, how can we remove redundancy and increase diversity in the resulting summarization during step two? The technique **maximal marginal relevance** (MMR) reranking can be applied to our problem. Essentially, this algorithm greedily reranks each sentence in the current passage, outputting only the top few as a summary. Figure 16.2 shows the output of the algorithm when we only select one sentence from each passage.

The MMR algorithm is as follows. Assume we are given an original list R and a profile p to construct the set of selected sentences S (where $|S| \ll |R|$). R is a partitioned chunk of sentences in the document we wish to summarize. The profile p determines what is exactly meant by “relevance.” Originally, the MMR formula was applied to documents returned from an information retrieval system (hence the term *reranking*). Documents were selected based on their marginal relevance to a query (which is our variable p) in addition to non-redundancy to already-selected documents. Since our task deals with sentence retrieval, p can be a user profile (text about the user), the entire document itself, or it could even be a query formulated by the user.

According to marginal relevance, the next sentence s_i to be added into the selected list S is defined as

$$s_i = \arg \max_{s \in R \setminus S} \left\{ (1 - \lambda) \cdot \text{sim}_1(s, p) - \lambda \cdot \arg \max_{s_j \in S} \text{sim}_2(s, s_j) \right\}. \quad (16.1)$$

The $R \setminus S$ notation may be read as “ R set minus S ”, i.e., all the elements in R that are not in S . The MMR formulation uses $\lambda \in [0, 1]$ to control relevance versus redundancy; the positive relevance score is discounted by the amount of redundancy (similarity) to the already-selected sentences. Again, the two similarity metrics may be any normalized, symmetric measures. The simplest instantiation for the similarity metric would be cosine similarity, and this is in fact the measure used in [Carbonell and Goldstein \[1998\]](#).

The algorithm may be terminated once an appropriate number of words or sentences is in S , or if the score $\text{sim}_1(s, p)$ is below some threshold. Furthermore, the similarity functions may be tweaked as well. Could you think of a way to include sentence position in the similarity function? That is, if a sentence is far away (dissimilar) from the candidate sentence, we could subtract from the similarity score. Even better, we could interpolate the two values into a new similarity score such as

$$\text{sim}(s, s') = \alpha \cdot \text{sim}_{\text{cosine}}(s, s') + (1 - \alpha) \cdot \left(1 - \frac{d(s, s')}{\max d(s, \cdot)} \right), \quad (16.2)$$

where $\alpha \in [0, 1]$ controls the weight between the regular cosine similarity and the distance measure, and $d(\cdot, \cdot)$ is the number of sentences between the two parameters. Note the “one minus” in front of the distance calculation, since a smaller distance implies a greater similarity.

Of course, λ in the MMR formula is also able to be set. In fact, for multi-document summarization, [Das and Martins \[2007\]](#) suggests starting out with $\lambda = 0.3$ and then slowly increasing to $\lambda = 0.7$. The reasoning behind this is to first emphasize novelty and then default to relevance. This should remind you of the exploration-exploitation tradeoff discussed in Chapter 11.

16.3 Abstractive Text Summarization

An abstractive summary creates sentences that did not exist in the original document or documents. Instead of a document vector, we will use a language model to represent the original text. Unlike the document vector, our language model gives us a principled way in which to generate text. Imagine we tokenized our document with unigram words. In our language model, we would have a parameter representing the probability of each word occurring. To create our own text, we will draw words from this probability distribution.

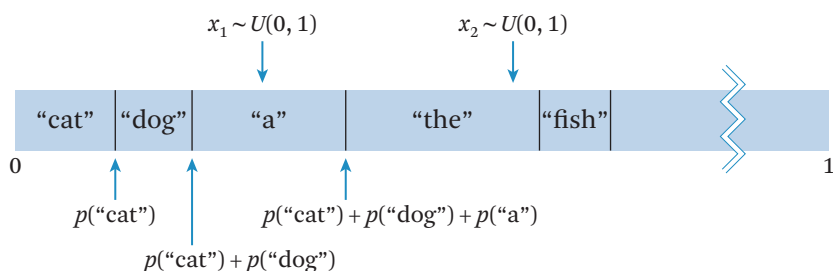


Figure 16.3 Drawing words from a unigram language model.

Say we have the unigram language model θ estimated on a document we wish to summarize. We wish to draw words w_1, w_2, w_3, \dots from θ that will comprise our summary. We want the word w_i to occur in our summary with about the same probability it occurred in the original document—this is how our generated text will approximate the longer document. Figure 16.3 depicts how we can accomplish this task. First, we create a list of all our parameters and incrementally sum their probabilities; this will allow us to use a random number on $[0, 1]$ to choose a word w_i . Simply, we get a uniform random floating point number between zero and one. Then, we iterate through the words in our vocabulary, summing their probabilities until we get to the random number. We output the term and repeat the process.

In the example, imagine we have the following values:

$p(\text{cat})$	0.010
$p(\text{cat}) + p(\text{dog})$	0.018
$p(\text{cat}) + p(\text{dog}) + p(\text{a})$	0.038
\vdots	\vdots
$p(\text{cat}) + p(\text{dog}) + p(\text{a}) + \dots + p(\text{zap})$	1.0

Say we generate a random number x_1 using a uniform distribution on $[0, 1]$. This is denoted as $x_1 \sim \mathcal{U}(0, 1)$. Now imagine that $x_1 = 0.032$. We go to the cumulative point 0.032 in our distribution and output “a”. We can repeat this process until our summary is of a certain length or until we generate an end-of-sentence token $\langle /s \rangle$.

At this point, you may be thinking that the text we generate will not make any sense—that is certainly true if we use a unigram language model since each word is generated independently without regard to its context. If more fluent language is required, we can use an n -gram language model, where $n > 1$. Instead of each word being independently generated, the new word will depend on the previous $n - 1$

words. The generation will work the same way as in the unigram case: say we have the word w_i and wish to generate w_{i+1} with a bigram language model. Our bigram language model gives us a distribution of words that occur after w_i and we draw the next word from there in the same way depicted in Figure 16.3.

The sentence generation from a bigram language model proceeds as follows: start with (e.g.) *The*. Then, pick from the distribution $p(w \mid \text{The})$ using the cumulative sum technique. The next selected word could be *cat*. Then, we use the distribution $p(w \mid \text{cat})$ to find the next w , and so on. While the unigram model only had one “sum table” (Figure 16.3), the bigram case needs V tables, one for each w' in $p(w \mid w')$.

Typically, the n -value will be around three to five depending on how much original data there is. We saw what happened when n is too small; we get a jumble of words that don’t make sense together. But we have another problem if n is too large. Consider the extreme case where $n = 20$. Then, given 19 words, we wish to generate the next one using our 20-gram language model. It’s very unlikely that those 19 words occurred more than once in our original document. That means there would only be one choice for the 20th word. Because of this, we would just be reproducing the original document, which is not a very good summary. In practice, we would like to choose an n -gram value that is large enough to produce coherent text yet small enough to not simply reproduce the corpus.

There is one major disadvantage to this abstractive summarization method. Due to its nature, a given word only depends on the n surrounding words. That is, there will be no long-range dependencies in our generated text. For example, consider the following sentence generated from a trigram language model:

They imposed a gradual smoking ban on virtually all corn seeds planted are hybrids.

All groups of three words make sense, but as a whole the sentence is incomprehensible; it seems the writer changed the topic from a smoking ban to hybrid crops mid-sentence. In special cases, when we restrict the length of a summary to a few words when summarizing highly redundant text, such a strategy appears to be effective as shown in the micropinion summarization method described in [Ganesan et al. \[2012\]](#).

16.3.1 Advanced Abstractive Methods

Some advanced abstractive methods rely more heavily on natural language processing to build a model of the document to summarize. **Named entity recognition** can be used to extract people, places, or businesses from the text. **Dependency parsers** and other syntactic techniques can be used to find the relation between the entities

and the actions they perform. Once these actors and roles are discovered, they are stored in some internal representation. To generate the actual text, some representations are chosen from the parsed collection, and English sentences are created based on them; this is called **realization**.

Such realization systems have much more fine-grained control over the generated text than the basic abstractive language model generator described above. A templated document structure may exist (such as intro→paragraph 1→paragraph 2→conclusion), and the structures are chosen to fill each spot. This control over text summarization and layout enables an easily-readable summary since it has a natural topical flow. In this environment, it would be possible to merge similar sentences with conjunctions such as *and* or *but*, depending on the context. To make the summary sound even more natural, pronouns can be used instead of entity names if the entity name has already been mentioned. Below are examples of these two operations:

Gold prices fell today. Silver prices fell today. → Gold and silver prices fell today.
Company A lost 9.43% today. Company A was the biggest mover. → Company A lost 9.43% today. It was the biggest mover.

Even better would be

Company A was today's biggest mover, losing 9.43%.

These operations are possible since the entities are stored in a structured format. For more on advanced natural language generation, we suggest [Reiter and Dale \[2000\]](#), which has a focus on practicality and implementation.

16.4 Evaluation of Text Summarization

In extractive summarization, representative sentences were selected from passages in the text and output as a summary. This solution is modeled as an information retrieval problem, and we can evaluate it as such. Redundancy is a critical issue, and the MMR technique we discussed attempts to alleviate it. When doing our evaluation, we should consider redundant sentences to be irrelevant, since the user does not want to read the same information twice. For a more detailed explanation of IR evaluation measures, please consult Chapter 9.

For full output scoring, we should prefer IR evaluation metrics that do not take into account result position. Although our summary is generated by ranked sentences per passage, the entire output is not a ranked list since the original document is composed of multiple passages. Therefore we can use precision, recall, and F_1 score.

It is possible to rank the passage scoring retrieval function using position-dependent metrics such as average precision or NDCG, but with the final output this is not feasible. Thus we need to decide whether to evaluate the passage scoring or the entire output (or both). Entire output scoring is likely more useful for actual users, while passage scoring could be useful for researchers to fine-tune their methods.

In abstractive summarization, we can't use the IR measures since we don't have a fixed set of candidate sentences. How can we compute recall if we don't know the total number of relevant sentences? There is also no intermediate ranking stage, so we also can't use average precision or NDCG (and again, we don't even know the complete set of correct sentences).

A laborious yet accurate evaluation would have human annotators create a gold standard summary. This "perfect" summary would be compared with the generated one, and some measure (e.g., ROUGE) would be used to quantify the difference. For the comparison measure, we have many possibilities—any measure that can compare two groups of text would be potentially applicable. For example, we can use the cosine similarity between the gold standard and generated summary. Of course, this has the downside that fluency is completely ignored (using unigram words). An alternative means would be to learn an n -gram language model over the gold standard summary, and then calculate the log-likelihood of the generated summary. This can ensure a basic level of fluency at the n -gram level, while also producing an interpretable result. Other comparisons between two probability distributions would also be applicable, such as KL-divergence.

The overall effectiveness of a summary can be tested if users read a summary and then answer questions about the original text. Was the summary able to capture the important information that the evaluator needs? If the original text was an entire textbook chapter, could the user read a three-paragraph summary and obtain sufficient information to answer the provided exercises? This is the only metric that can be used for both extractive and abstractive measures. Using a language model to score an extractive summary vs. an abstractive one would likely be biased towards the extractive one since this method contains phrases directly from the original text, giving it a very high likelihood.

16.5 Applications of Text Summarization

At the beginning of the chapter, we've already touched on a few summarization applications; we mentioned news articles, retrieval results, and opinion summarization. Summarization saves users time from manually reading the entire corpus while simultaneously enhancing preexisting data with summary "annotations."

The aspect opinion analysis mentioned earlier segments portions of user reviews into speaking about a particular topic. We can use this topic analysis to collect passages of text into a large group of comments on one aspect. Instead of describing this aspect with sorted unigram words, we could run a summarizer on each topic, generating readable text as output. These two methods complement each other, since the first step finds what aspects the users are interested in, while the second step conveys the information.

A theme in this book is the union of both structured and unstructured data, mentioned much more in detail in Chapter 19. Summarization is an excellent example of this application. For example, consider a financial summarizer with text reports from the Securities and Exchange Commission (SEC) as well as raw stock market data. Summarizing both these data sources in one location would be very valuable for (e.g.) mutual fund managers or other financial workers. Being able to summarize (in text) a huge amount of structured trading data could reveal patterns that humans would otherwise be unaware of—this is an example of **knowledge discovery**.

E-discovery (electronic discovery) is the process of finding relevant information in litigation (lawsuits and court cases). Lawyers rely on e-discovery to sift through vast amounts of textual information to build their case. The Enron email dataset¹ is a well-known corpus in this field. Summarizing email correspondence between two people or a department lets investigators quickly decide whether they'd like to dig deeper in a particular area or try another approach. In this way, summarization and search are coupled; search allows a subset of data to be selected that is relevant to a query, and the summarization can take the search results and quickly explain them to the user. Finally, linking email correspondence together (from sender to receivers) is a structured complement to the unstructured text content of the email itself.

Perhaps of more interest to those reading this book is the ability to summarize research from a given field. Given proceedings from a conference, could we have a summarizer explain the main trends and common approaches? What was most novel compared to previous conferences? When writing your own paper, can you write everything except the introduction and related work? The introduction is an overview summary of your paper. Related work is mostly a summary of papers similar to yours.

1. <https://www.cs.cmu.edu/~./enron/>

Bibliographic Notes and Further Reading

As mentioned in this chapter, [Das and Martins \[2007\]](#) is a comprehensive survey on summarization techniques. Additionally, [Nenkova and McKeown \[2012\]](#) is a valuable read. For applications, latent aspect rating analysis [[Wang et al. 2010](#)], [[Wang et al. 2011](#)] is a form of summarization applied to product reviews. We mention this particular application in more detail in Chapter 18. A typical extractive summarizer is presented in [Radev et al. \[2004\]](#), a typical abstractive summarizer is presented in [Ganesan et al. \[2010\]](#), and evaluation suggestions are presented in [Steinberger and Jezek \[2009\]](#). The MMR algorithm was originally described in [Carbonell and Goldstein \[1998\]](#). For advanced NLG (natural language generation) techniques, a good starting point is [Reiter and Dale \[2000\]](#).

Exercises

- 16.1. Do you think one summarization method (extractive or abstractive) would perform better on a small dataset? How about a large dataset? Justify your reasoning.
- 16.2. Explain how you can improve the passage detection by looking beyond only the adjacent sentences. How would you implement this?
- 16.3. Write a basic passage segmenter in META. As input, take a document and extract the sentences into a vector with a built-in tokenizer. Segment the vector into passages using a similarity algorithm.
- 16.4. Now that you have a document segmented into passages, use META to set up a search engine over each passage, where you treat passages as individual documents. Ensure that you have enough sentences per passage. You may need to tweak your previous answer to achieve this.
- 16.5. With your passage search engine, find a representative sentence from each passage to create a summary for the original document.
- 16.6. Use META's language model to learn a distribution of words over a document you wish to summarize.
- 16.7. Add a `generate` function to the language model. It should take a context ($n - 1$ terms) and generate the n^{th} term. Use the calculation described in this chapter to generate the next word.
- 16.8. Summarize the input document using the generator. Experiment with different stopping criteria. Which seems to work the best?

16.9. Create some simple post-processing rules for natural language generation realization. The examples we gave in the text were sentence joining and pronoun insertion. What else can you think of?

16.10. Explain how we can combine text summarization and topic modeling to create a powerful exploratory text mining application.

16.11. What can we accomplish by interpolating a language model distribution for an abstractive summarizer with another probability distribution, perhaps from existing summaries?

References

- C. C. Aggarwal. 2015. *Data Mining - The Textbook*. Springer. DOI: [10.1007/978-3-319-14142-8](https://doi.org/10.1007/978-3-319-14142-8). 296
- C. C. Aggarwal and C. Zhai, editors. 2012. *Mining Text Data*. Springer. DOI: [10.1007/978-1-4614-3223-4](https://doi.org/10.1007/978-1-4614-3223-4). 296, 315
- J. Allen. 1995. *Natural Language Understanding*. 2nd ed. Benjamin-Cummings Publishing Co., Inc., Redwood City, CA. 54
- G. Amati and C. J. Van Rijsbergen. October 2002. Probabilistic models of information retrieval based on measuring the divergence from randomness. *ACM Trans. Inf. Syst.*, 20(4):357–389. DOI: [10.1145/582415.582416](https://doi.org/10.1145/582415.582416). 87, 88, 90, 111
- A. U. Asuncion, M. Welling, P. Smyth, and Y. W. Teh. 2009. On smoothing and inference for topic models. In *UAI 2009, Proc. of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence, Montreal, QC, Canada, June 18-21, 2009*, pp. 27–34. 385
- R. A. Baeza-Yates and B. A. Ribeiro-Neto. 2011. *Modern Information Retrieval - the concepts and technology behind search*. 2nd ed. Pearson Education Ltd., Harlow, UK. <http://www.mir2ed.org/>. xvii, 18, 19
- Y. Bar-Hillel, *The Present Status of Automatic Translation of Languages*, in *Advances in Computers*, vol. 1 (1960), pp. 91–163.
- R. Belew. 2008. *Finding Out About: A Cognitive Perspective on Search Engine Technology and the WWW*. Cambridge University Press. 18
- N. J. Belkin and W. B. Croft. 1992. Information filtering and information retrieval: Two sides of the same coin? *Commun. ACM*, 35(12):29–38. DOI: [10.1145/138859.138861](https://doi.org/10.1145/138859.138861). 84
- C. M. Bishop. 2006. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag New York, Inc., Secaucus, NJ. 19, 37, 312, 385, 462
- D. M. Blei, A. Y. Ng, and M. I. Jordan. March 2003. Latent Dirichlet Allocation. *J. of Mach. Learn. Res.*, 3:993–1022. 385

- J. S. Breese, D. Heckerman, and C. Kadie. 1998. Empirical analysis of predictive algorithms for collaborative filtering. In *Proc. of the Fourteenth Conference on Uncertainty in Artificial Intelligence*, UAI'98, Morgan Kaufmann Publishers Inc. pp. 43–52, San Francisco, CA. <http://dl.acm.org/citation.cfm?id=2074094.2074100>. 235
- P. F. Brown, P. V. deSouza, R. L. Mercer, V. J. Della Pietra, and J. C. Lai. 1992. Class-based N-gram Models of Natural Language. *Comput. Linguist.*, 18(4):467–479. 273, 288, 290, 291
- C. Buckley. 1994. Automatic query expansion using smart: Trec 3. In *Proc. of The third Text REtrieval Conference (TREC-3)*, pp. 69–80. 144
- S. Büttcher, C. Clarke, and G. V. Cormack. 2010. *Information Retrieval: Implementing and Evaluating Search Engines*. The MIT Press. xvii, 18, 165
- F. Cacheda, V. Carneiro, D. Fernández, and V. Formoso. 2011. Comparison of collaborative filtering algorithms: Limitations of current techniques and proposals for scalable, high-performance recommender systems. *ACM Trans. Web*, 5(1):2:1–2:33. DOI: [10.1145/1921591.1921593](https://doi.org/10.1145/1921591.1921593). 235
- C. Campbell and Y. Ying. 2011. *Learning with Support Vector Machines*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers. DOI: [10.2200/S00324ED1V01Y201102AIM010](https://doi.org/10.2200/S00324ED1V01Y201102AIM010). 311
- J. Carbonell and J. Goldstein. 1998. The Use of MMR, Diversity-based Reranking for Reordering Documents and Producing Summaries. In *Proc. of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '98, ACM, pp. 335–336, New York. DOI: [doi=10.1.1.188.3982](https://doi.org/10.1.1.188.3982) 321, 327
- C.-C. Chang and C.-J. Lin. 2011. LIBSVM: A Library for Support Vector Machines. *ACM Trans. Intell. Syst. Technol.*, 2(3):27:1–27:27. 58
- J. Chang, S. Gerrish, C. Wang, J. L. Boyd-graber, and D. M. Blei. 2009. Reading Tea Leaves: How Humans Interpret Topic Models. In Y. Bengio, D. Schuurmans, J.D. Lafferty, C.K.I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, Curran Associates, Inc. 22, pp. 288–296. 272, 383, 384, 385, 410
- K. W. Church and P. Hanks. 1990. Word association norms, mutual information, and lexicography. *Comput. Linguist.*, 16(1):22–29. <http://dl.acm.org/citation.cfm?id=89086.89095>. 273
- T. Cover and J. Thomas. 1991. *Elements of Information Theory*. New York: Wiley. DOI: [DOI: 10.1002/047174882X](https://doi.org/10.1002/047174882X) 37, 473
- B. Croft, D. Metzler, and T. Strohman. 2009. *Search Engines: Information Retrieval in Practice*, 1st ed., Addison-Wesley Publishing Company. xvii, 18, 165

- D. Das and A. F. T. Martins. 2007. A Survey on Automatic Text Summarization. Technical report, Literature Survey for the Language and Statistics II course at Carnegie Mellon University. [318](#), [321](#), [327](#)
- R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. 2008. LIBLINEAR: A Library for Large Linear Classification. *J. Mach. Learn. Res.*, 9:1871–1874. [58](#)
- H. Fang, T. Tao, and C. Zhai. 2004. A formal study of information retrieval heuristics. In *Proc. of the 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '04, ACM, pp. 49–56, New York. DOI: [10.1145/1008992.1009004](#). [129](#)
- H. Fang, T. Tao, and C. Zhai. April 2011. Diagnostic evaluation of information retrieval models. *ACM Trans. Inf. Syst.*, 29(2):7:1–7:42. DOI: [10.1145/1961209.1961210](#). [88](#), [90](#), [129](#)
- R. Feldman and J. Sanger. 2007. *The Text Mining Handbook - Advanced Approaches in Analyzing Unstructured Data*. Cambridge University Press. [18](#)
- E. A. Fox, M. A. Gonçalves, and R. Shen. 2012. *Theoretical Foundations for Digital Libraries: The 5S (Societies, Scenarios, Spaces, Structures, Streams) Approach*. Synthesis Lectures on Information Concepts, Retrieval, and Services. Morgan & Claypool Publishers. DOI: [10.2200/S00434ED1V01Y201207ICR022](#). [80](#)
- W. B. Frakes and R. A. Baeza-Yates, editors. 1992. *Information Retrieval: Data Structures & Algorithms*. Prentice-Hall, [18](#)
- K. Ganesan, C. Zhai, and J. Han. 2010. Opinions: A graph-based approach to abstractive summarization of highly redundant opinions. In *Proc. of the 23rd International Conference on Computational Linguistics*, COLING '10, Association for Computational Linguistics, pp. 340–348, Stroudsburg, PA. [327](#)
- K. Ganesan, C. Zhai, and E. Viegas. 2012. Micropinion generation: an unsupervised approach to generating ultra-concise summaries of opinions. In *Proc. of the 21st World Wide Web Conference 2012, WWW 2012*, Lyon, France, April 16–20, 2012, pages 869–878. DOI: [10.1145/2187836.2187954](#) [323](#)
- J. Gantz, and D. Reinsel. 2012. The Digital Universe in 2020: Big Data, Bigger Digital Shadows, and Biggest Growth in the Far East, IDC Report, December, 2012. [3](#)
- A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. 1995. *Bayesian Data Analysis*. Chapman & Hall. [37](#)
- S. Ghemawat, H. Gobioff, and S.-T. Leung. 2003. *The Google file system*. In *Proc. of the nineteenth ACM symposium on Operating systems principles* (SOSP '03). ACM, New York, 29–43. [195](#)
- M. A. Gonçalves, E. A. Fox, L. T. Watson, and N. A. Kipp. 2004. Streams, structures, spaces, scenarios, societies (5s): A formal model for digital libraries. *ACM Trans. Inf. Syst.*, 22(2):270–312. DOI: [10.1145/984321.984325](#). [84](#)

- D. A. Grossman and O. Frieder. Kluwer, 2004. *Information Retrieval - Algorithms and Heuristics, Second Edition*, vol. 15 of *The Kluwer International Series on Information Retrieval*. DOI: [10.1007/978-1-4020-3005-5](https://doi.org/10.1007/978-1-4020-3005-5). 18
- G. Hamerly and C. Elkan. 2003. Learning the k in k-means. In *Advances in Neural Information Processing Systems 16 [Neural Information Processing Systems, NIPS December 8-13, 2003]*, Vancouver and Whistler, British Columbia, Canada], pp. 281–288. DOI: [doi=10.1.1.9.3574](https://doi.org/10.1.1.9.3574) 295
- J. Han. 2005. *Data Mining: Concepts and Techniques*. Morgan Kaufmann Publishers Inc., San Francisco, CA. 296
- D. Harman. 2011. *Information Retrieval Evaluation*. Synthesis Lectures on Information Concepts, Retrieval, and Services. Morgan & Claypool Publishers. DOI: [10.1145/215206.215351](https://doi.org/10.1145/215206.215351) 168, 188
- M. A. Hearst. 2009. *Search User Interfaces*. 1st ed. Cambridge University Press, New York. 19, 85
- J. L. Herlocker, J. A. Konstan, L. G. Terveen, and J. T. Riedl. 2004. Evaluating Collaborative Filtering Recommender Systems. *ACM Trans. Inf. Syst.*, 22(1):5–53. DOI: [10.1145/963770.963772](https://doi.org/10.1145/963770.963772) 235
- J. L. Hodges and E. L. Lehmann. 1970. *Basic Concepts of Probability and Statistics*. Holden Day, San Francisco. 36
- T. Hofmann. 1999. Probabilistic Latent Semantic Analysis. In *Proc. of the Fifteenth Conference on Uncertainty in Artificial Intelligence, UAI'99*, Morgan Kaufmann Publishers Inc., pp. 289–296, San Francisco, CA. DOI: [10.1145/312624.312649](https://doi.org/10.1145/312624.312649) 370, 385
- A. Huang. 2008. Similarity Measures for Text Document Clustering. In *Proc. of the Sixth New Zealand Computer Science Research Student Conference (NZCSRSC2008)*, Christchurch, New Zealand, pages 49–56. 280
- F. Jelinek. 1997. *Statistical Methods for Speech Recognition*. MIT Press, Cambridge, MA. 30, 54
- J. Jiang. 2012. Information extraction from text, In Charu C. Aggarwal and ChengXiang Zhai (Eds.), *Mining Text Data*, Springer, pp. 11–41. 19, 55
- S. Jiang and C. Zhai. 2014. Random walks on adjacency graphs for mining lexical relations from big text data. In *2014 IEEE International Conference on Big Data, Big Data 2014*, Washington, DC, USA, October 27-30, pages 549–554. DOI: [10.1109/BigData.2014.7004272](https://doi.org/10.1109/BigData.2014.7004272). 273
- Y. Jo and A. H. Oh. 2011. Aspect and sentiment unification model for online review analysis. In *Proceedings of the Fourth ACM International Conference on Web Search and Data Mining, WSDM '11*, ACM, pp. 815–824, New York. DOI: [10.1145/1935826.1935932](https://doi.org/10.1145/1935826.1935932). 410

- T. Joachims, L. Granka, B. Pan, H. Hembrooke, F. Radlinski, and G. Gay. 2007. Evaluating the accuracy of implicit feedback from clicks and query reformulations in web search. *ACM Trans. Inf. Syst.*, 25(2). DOI: [10.1145/1229179.1229181](https://doi.org/10.1145/1229179.1229181). 144
- D. Jurafsky and J. H. Martin. 2009. *Speech and Language Processing*. 2nd ed. Prentice-Hall, Inc., Upper Saddle River, NJ. 19, 54
- D. Kelly. 2009. Methods for Evaluating Interactive Information Retrieval Systems with Users. *Foundations and Trends in Information Retrieval*, 3(1-2):1–224. DOI: [10.1561/15000000012](https://doi.org/10.1561/15000000012) 168, 188
- D. Kelly and J. Teevan. 2003. Implicit feedback for inferring user preference: A bibliography. *SIGIR Forum*, 37(2):18–28. DOI: [10.1145/959258.959260](https://doi.org/10.1145/959258.959260). 144
- H. D. Kim, M. Castellanos, M. Hsu, C. Zhai, T. Rietz, and D. Diermeier. 2013. Mining causal topics in text data: iterative topic modeling with time series feedback. In *Proc. of the 22nd ACM international conference on Conference on information and knowledge management*, CIKM '13, ACM pages 885–890, New York, NY. DOI: [10.1145/2505515.2505612](https://doi.org/10.1145/2505515.2505612). 435, 438, 439, 440
- J. M. Kleinberg. 1999. Authoritative sources in a hyperlinked environment. *J. ACM*, 46(5):604–632. DOI: [10.1145/324133.324140](https://doi.org/10.1145/324133.324140). 216
- J. M. Kleinberg. 2002. An impossibility theorem for clustering. In *Advances in Neural Information Processing Systems 15 [Neural Information Processing Systems, NIPS 2002, December 9-14, 2002, Vancouver, British Columbia, Canada]*, pp. 446–453. <http://papers.nips.cc/paper/2340-an-impossibility-theorem-for-clustering>. 296
- D. Koller and N. Friedman. 2009. *Probabilistic Graphical Models: Principles and Techniques - Adaptive Computation and Machine Learning*. The MIT Press. 385
- J. Lafferty and C. Zhai. 2003. Probabilistic relevance models based on document and query generation. In W. Bruce Croft and John Lafferty, editors, *Language Modeling and Information Retrieval*. Kluwer Academic Publishers. DOI: [10.1007/978-94-017-0171-6_1](https://doi.org/10.1007/978-94-017-0171-6_1) 87, 113
- D. Lin. 1999. Automatic identification of non-compositional phrases. In *Proc. of the 37th Annual Meeting of the Association for Computational Linguistics on Computational Linguistics*, ACL '99, Association for Computational Linguistics, pages 317–324, Stroudsburg, PA. DOI: [10.3115/1034678.1034730](https://doi.org/10.3115/1034678.1034730). 273, 291
- J. Lin and C. Dyer. 2010. *Data-Intensive Text Processing with MapReduce*. Morgan and Claypool Publishers. DOI: [10.2200/S00274ED1V01Y201006HLT007](https://doi.org/10.2200/S00274ED1V01Y201006HLT007). 198, 216
- Bing Liu. 2012. *Sentiment Analysis and Opinion Mining*. Synthesis Lectures on Human Language Technologies. Morgan & Claypool Publishers. DOI: [10.2200/S00416ED1V01Y201204HLT016](https://doi.org/10.2200/S00416ED1V01Y201204HLT016). 410
- T.-Y. Liu. 2009. Learning to rank for information retrieval. *Found. Trends Inf. Retr.*, 3(3):225–331. DOI: [10.1561/15000000016](https://doi.org/10.1561/15000000016). 216

- Y. Lv and C. Zhai. 2009. A comparative study of methods for estimating query language models with pseudo feedback. In *Proc. of the 18th ACM Conference on Information and Knowledge Management*, CIKM '09, ACM, pp. 1895–1898, New York. DOI: [10.1145/1645953.1646259](https://doi.org/10.1145/1645953.1646259). 144
- Y. Lv and C. Zhai. 2010. Positional relevance model for pseudo-relevance feedback. In *Proc. of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '10, ACM, pages 579–586, New York. DOI: [10.1145/1835449.1835546](https://doi.org/10.1145/1835449.1835546). 144
- Y. Lv and C. Zhai. 2011. Lower-bounding Term Frequency Normalization. In *Proc. of the 20th ACM International Conference on Information and Knowledge Management*, CIKM '11, pp. 7–16. DOI: [10.1145/2063576.2063584](https://doi.org/10.1145/2063576.2063584) 88, 110
- P. Lyman, H. R. Varian, K. Swearingen, P. Charles, N. Good, L.L. Jordan, and J. Pal. 2003. How much information? <http://www2.sims.berkeley.edu/research/projects/how-much-info-2003>. 3
- C. D. Manning and H. Schütze. 1999. *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge, MA. 19, 54, 273
- C. D. Manning, P. Raghavan, and H. Schütze. 2008. *Introduction to Information Retrieval*. Cambridge University Press, New York. xvii, 18, 165, 315
- M. E. Maron and J. L. Kuhns. 1960. On relevance, probabilistic indexing and information retrieval. *Journal of the ACM*, 7:216–244. DOI: [10.1145/321033.321035](https://doi.org/10.1145/321033.321035) 87
- S. Massung and C. Zhai. 2015. SyntacticDiff: Operator-Based Transformation for Comparative Text Mining. In *Proc. of the 3rd IEEE International Conference on Big Data*, pp. 571–580. 306
- S. Massung and C. Zhai. 2016. Non-Native Text Analysis: A Survey. *The Journal of Natural Language Engineering*, 22(2):163–186. DOI: [10.1017/S1351324915000303](https://doi.org/10.1017/S1351324915000303) 306
- S. Massung, C. Zhai, and J. Hockenmaier. 2013. Structural Parse Tree Features for Text Representation. In *IEEE Seventh International Conference on Semantic Computing*, pp. 9–13. DOI: [10.1109/ICSC.2013.13](https://doi.org/10.1109/ICSC.2013.13) 305
- J. D. McAuliffe and D. M. Blei. 2008. Supervised topic models. In J.C. Platt, D. Koller, Y. Singer, and S.T. Roweis, eds., *Advances in Neural Information Processing Systems 20*, pages 121–128. Curran Associates, Inc. 386
- G. J. McLachlan and T. Krishnan. 2008. *The EM algorithm and extensions*. 2nd ed. Wiley Series in Probability and Statistics. Hoboken, NJ., Wiley. http://gso.gbv.de/DB=2.1/CMD?ACT=SRCHA&SRT=YOP&IKT=1016&TRM=ppn+52983362X&sourceid=fbw_bibsonomy. DOI: [10.1002/9780470191613](https://doi.org/10.1002/9780470191613) 466
- Q. Mei. 2009. Contextual text mining. Ph.D. Dissertation, University of Illinois at Urbana-Champaign. 440

- Q. Mei and C. Zhai. 2006. A mixture model for contextual text mining. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '06, ACM, pp. 649–655, New York. DOI: [10.1145/1150402.1150482](https://doi.org/10.1145/1150402.1150482). 423, 440
- Q. Mei, D. Xin, H. Cheng, J. Han, and C. Zhai. 2006. Generating semantic annotations for frequent patterns with context analysis. In *Proc. of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '06, ACM, pp. 337–346, New York. DOI: [10.1145/1150402.1150441](https://doi.org/10.1145/1150402.1150441). 417
- Q. Mei, C. Liu, H. Su, and C. Zhai. 2006. A probabilistic approach to spatiotemporal theme pattern mining on weblogs. In *Proc. of the 15th international conference on World Wide Web (WWW '06)*. ACM, New York, 533–542. DOI: [10.1145/1135777.1135857](https://doi.org/10.1145/1135777.1135857). 425, 426
- Q. Mei, X. Ling, M. Wondra, H. Su, and C. Zhai. 2007a. Topic sentiment mixture: Modeling facets and opinions in weblogs. In *Proc. of the 16th International Conference on World Wide Web*, WWW '07, ACM, pp. 171–180, New York. DOI: [10.1145/1242572.1242596](https://doi.org/10.1145/1242572.1242596). 410
- Q. Mei, X. Shen, and C. Zhai. 2007b. Automatic labeling of multinomial topic models. In *Proc. of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, San Jose, California, August 12–15, 2007, pp. 490–499. DOI: [10.1145/1281192.1281246](https://doi.org/10.1145/1281192.1281246). 278
- Q. Mei, D. Cai, D. Zhang, and C. Zhai. 2008. Topic modeling with network regularization. In *Proceedings of the 17th International Conference on World Wide Web*, WWW '08, ACM, pp. 101–110, New York. DOI: [10.1145/1367497.1367512](https://doi.org/10.1145/1367497.1367512). 431, 432, 440
- T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur. 2010. Recurrent neural network based language model. In *INTERSPEECH 2010, 11th Annual Conference of the International Speech Communication Association*, Makuhari, Chiba, Japan, September 26–30, 2010, pp. 1045–1048. http://www.isca-speech.org/archive/interspeech_2010/i10_1045.html. 292
- T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. 2013. Distributed Representations of Words and Phrases and their Compositionality. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013*. Proceedings of a meeting held December 5–8, 2013, Lake Tahoe, NV, pp. 3111–3119. 273, 292, 293
- T. M. Mitchell. 1997. *Machine learning*. McGraw Hill Series in Computer Science. McGraw-Hill. 19, 37, 315
- M.-F. Moens. 2006. *Information Extraction: Algorithms and Prospects in a Retrieval Context (The Information Retrieval Series)*. Springer-Verlag New York, Inc., Secaucus, NJ. DOI: [10.1007/978-1-4020-4993-4](https://doi.org/10.1007/978-1-4020-4993-4). 55

- I. J. Myung. 2003. Tutorial on maximum likelihood estimation. *J. Math. Psychol.*, 47(1):90–100. DOI: [10.1016/S0022-2496\(02\)00028-7](https://doi.org/10.1016/S0022-2496(02)00028-7). 36
- A. Nenkova and K. McKeown. 2012. A survey of text summarization techniques. In Charu C. Aggarwal and C. Zhai, eds, *Mining Text Data*, pp. 43–76. Springer US. DOI: [10.1007/978-1-4614-3223-4_3](https://doi.org/10.1007/978-1-4614-3223-4_3). 327
- L. Page, S. Brin, R. Motwani, and T. Winograd. 1999. The PageRank Citation Ranking: Bringing Order to the Web. <http://ilpubs.stanford.edu:8090/422/1/1999-66.pdf>. 216
- B. Pang and L. Lee. 2008. Opinion Mining and Sentiment Analysis. *Foundations and Trends in Information Retrieval*, 2(1-2):1–135. DOI: [10.1561/15000000011](https://doi.org/10.1561/15000000011). 409, 410
- J. M. Ponte and W. B. Croft. 1998. A language modeling approach to information retrieval. In *Proc. of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '98, ACM, pp. 275–281, New York, NY. DOI: [10.1145/290941.291008](https://doi.org/10.1145/290941.291008). 87, 90, 128, 427
- J. R. Quinlan. 1986. Induction of Decision Trees. *Machine Learning*, 1(1):81–106. DOI: [10.1007/BF00116251](https://doi.org/10.1007/BF00116251). 301
- D. R. Radev, H. Jing, M. Styś, and D. Tam. 2004. Centroid-based summarization of multiple documents. *Information Processing & Management*, 40(6):919–938. DOI: [10.1016/j.ipm.2003.10.006](https://doi.org/10.1016/j.ipm.2003.10.006). 327
- D. Ramage, D. Hall, R. Nallapati, and C. D. Manning. 2009. Labeled lda: A supervised topic model for credit attribution in multi-labeled corpora. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 1 - Volume 1*, EMNLP '09, Association for Computational Linguistics, pages 248–256, Stroudsburg, PA. 386
- E. Reiter and R. Dale. 2000. *Building Natural Language Generation Systems*. Cambridge University Press, New York. 324, 327
- F. Ricci, L. Rokach, B. Shapira, and P. B. Kantor. 2010. *Recommender Systems Handbook*. 1st ed. Springer-Verlag New York, Inc. DOI: [10.1007/978-0-387-85820-3](https://doi.org/10.1007/978-0-387-85820-3). 235
- C. J. Van Rijsbergen. 1979. *Information Retrieval*. 2nd ed. Butterworth-Heinemann, Newton, MA.
- S. Robertson and K. Sparck Jones. 1976. Relevance weighting of search terms. *Journal of the American Society for Information Science*, 27:129–146. 87
- S. E. Robertson. 1997. Readings in Information Retrieval. In *The Probability Ranking Principle in IR*, San Francisco, CA, Morgan Kaufmann Publishers Inc. pp. 281–286. 84, 85
- S. Robertson and H. Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Found. Trends Inf. Retr.*, 3(4):333–389. DOI: [10.1561/15000000019](https://doi.org/10.1561/15000000019). 88, 89, 129

- S. Robertson, H. Zaragoza, and M. Taylor. 2004. Simple BM25 Extension to Multiple Weighted Fields. In *Proc. of the Thirteenth ACM International Conference on Information and Knowledge Management*, CIKM '04, pp. 42–49. DOI: [10.1145/1031171.1031181](https://doi.org/10.1145/1031171.1031181) 110
- C. Roe. 2012. The growth of unstructured data: what to do with all those zettabytes? <http://www.dataversity.net/the-growth-of-unstructured-data-what-are-we-going-to-do-with-all-those-zettabytes/>. 3
- R. Rosenfeld. 2000. Two decades of statistical language modeling: Where do we go from here. In *Proceedings of the IEEE*. 54
- G. Salton. 1989. *Automatic Text Processing: The Transformation, Analysis and Retrieval of Information by Computer*. Addison-Wesley. 18
- G. Salton and M. McGill. 1983. *Introduction to Modern Information Retrieval*. McGraw-Hill. 18
- G. Salton, A. Wong, and C. S. Yang. 1975. A vector space model for automatic indexing. *Commun. ACM*, 18(11):613–620. 87
- G. Salton and C. Buckley. 1990. Improving retrieval performance by relevance feedback. *Journal of the American Society for Information Science*, 41:288–297. 144
- M. Sanderson. 2010. Test Collection Based Evaluation of Information Retrieval Systems. *Foundations and Trends in Information Retrieval*, 4(4):247–375. 168, 188
- M. Sanderson and W. B. Croft. 2012. The history of information retrieval research. *Proc. of the IEEE*, 100(Centennial-Issue):1444–1451, 2012. DOI: [10.1109/JPROC.2012.2189916](https://doi.org/10.1109/JPROC.2012.2189916). 85
- S. Sarawagi. 2008. Information extraction. *Found. Trends databases*, 1(3):261–377. DOI: [10.1561/19000000003](https://doi.org/10.1561/19000000003). 19, 55
- F. Sebastiani. 2002. Machine learning in automated text categorization. *ACM Comput. Surv.*, 34(1):1–47. DOI: [10.1145/505282.505283](https://doi.org/10.1145/505282.505283). 315
- G. Shani and A. Gunawardana. 2011. Evaluating Recommendation Systems. In *Recommender Systems Handbook*, 2nd ed., pp. 257–297. Springer, New York, NY. DOI: [10.1007/978-0-387-85820-3_8](https://doi.org/10.1007/978-0-387-85820-3_8). 235
- F. Silvestri. 2010. Mining query logs: Turning search usage data into knowledge. *Found. Trends Inf. Retr.*, 4:1–174. DOI: [10.1561/15000000013](https://doi.org/10.1561/15000000013) 144
- A. Singhal, C. Buckley, and Mandar Mitra. 1996. Pivoted document length normalization. In *Proc. of the 19th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '96, ACM, pp. 21–29, New York. DOI: [10.1145/243199.243206](https://doi.org/10.1145/243199.243206). 89, 106, 128
- N. Smith. 2010. Text-driven forecasting. <http://www.cs.cmu.edu/~ñasmith/papers/smith.whitepaper10.pdf>. 440

- Mark D. Smucker, James Allan, and Ben Carterette. 2007. A Comparison of Statistical Significance Tests for Information Retrieval Evaluation. In *Proc. of the Sixteenth ACM Conference on Conference on Information and Knowledge Management*, CIKM '07, ACM, pp. 623–632, New York. DOI: [10.1145/1321440.1321528](https://doi.org/10.1145/1321440.1321528). 185
- K. Sparck Jones and P. Willett, eds. 1997. *Readings in Information Retrieval*. San Francisco, CA, Morgan Kaufmann Publishers Inc. 18, 85, 188
- N. Spirin and J. Han. May 2012. Survey on Web Spam Detection: Principles and Algorithms. *SIGKDD Explor. Newsl.*, 13(2):50–64. DOI: [10.1145/2207243.2207252](https://doi.org/10.1145/2207243.2207252). 191
- E. Stamatatos. 2009. A Survey of Modern Authorship Attribution Methods. *J. Am. Soc. Inf. Sci. Technol.*, 60(3):538–556. DOI: [10.1002/asi.v60:3](https://doi.org/10.1002/asi.v60:3) 305
- M. Steinbach, G. Karypis, and V. Kumar. 2000. A comparison of document clustering techniques. In *KDD Workshop on Text Mining*. 296
- J. Steinberger and K. Jezek. 2009. Evaluation measures for text summarization. *Computing and Informatics*, 28(2):251–275. 327
- M. Steyvers and T. Griffiths. 2007. Probabilistic topic models. *Handbook of Latent Semantic Analysis*, 427(7):424–440. 385, 386
- Y. Sun and J. Han. 2012. *Mining Heterogeneous Information Networks: Principles and Methodologies*. Morgan & Claypool Publishers. DOI: [10.2200/S00433ED1V01Y201207DMK005](https://doi.org/10.2200/S00433ED1V01Y201207DMK005). 440
- I. Titov and R. McDonald. 2008. Modeling online reviews with multi-grain topic models. In *Proc. of the 17th International Conference on World Wide Web*, WWW '08, ACM, pp. 111–120, New York. DOI: [10.1145/1367497.1367513](https://doi.org/10.1145/1367497.1367513). 410
- H. Turtle and W. B. Croft. 1990. Inference networks for document retrieval. In *Proc. of the 13th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '90, ACM, pp. 1–24, New York. DOI: [10.1145/96749.98006](https://doi.org/10.1145/96749.98006). 88
- Princeton University. 2010. About wordnet. <http://wordnet.princeton.edu>. 395
- C. J. van Rijsbergen. 1979. *Information Retrieval*. Butterworths. 18
- H. Wang, Yue Lu, and C. Zhai. 2010. Latent Aspect Rating Analysis on Review Text Data: A Rating Regression Approach. In *Proc. of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '10, ACM, pp. 783–792, New York. DOI: [10.1145/1835804.1835903](https://doi.org/10.1145/1835804.1835903). 318, 327, 405, 406, 407, 408, 409, 410
- H. Wang, Y. Lu, and C. Zhai. 2011. Latent Aspect Rating Analysis Without Aspect Keyword Supervision. In *Proc. of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '11, ACM, pp. 618–626, New York. DOI: [10.1145/2020408.2020505](https://doi.org/10.1145/2020408.2020505). 318, 327, 405, 410

- J. Weizenbaum. 1966. ELIZA—A Computer Program for the Study of Natural Language Communication Between Man and Machine, *Communications of the ACM* 9 (1): 36–45, DOI: [10.1145/265153.365168](https://doi.org/10.1145/265153.365168). 44
- J. S. Whissell and C. L. A. Clarke. 2013. Effective Measures for Inter-document Similarity. In *Proc. of the 22nd ACM International Conference on Conference on Information & Knowledge Management, CIKM '13*, ACM, pages 1361–1370, New York. DOI: [10.1145/2505515.2505526](https://doi.org/10.1145/2505515.2505526). 279
- R. W. White and R. A. Roth. 2009. *Exploratory Search: Beyond the Query-Response Paradigm*. Synthesis Lectures on Information Concepts, Retrieval, and Services. Morgan & Claypool Publishers. DOI: [10.2200/S00174ED1V01Y200901ICR003](https://doi.org/10.2200/S00174ED1V01Y200901ICR003). 85
- R. W. White, B. Kules, S. M. Drucker, and m.c. schraefel. 2006. Introduction. *Commun. ACM*, 49(4):36–39. DOI: [10.1145/1121949.1121978](https://doi.org/10.1145/1121949.1121978). 85
- I. H. Witten, A. Moffat, and T. C. Bell. 1999. *Managing Gigabytes (2Nd Ed.): Compressing and Indexing Documents and Images*. Morgan Kaufmann Publishers Inc., San Francisco, CA. 18, 165
- C.F J. Wu. 1983. On the convergence properties of the EM algorithm. *Ann. of stat.*, 95–103. 368
- J. Xu and W. B. Croft. 1996. Query expansion using local and global document analysis. In *Proc. of the 19th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '96, ACM, pp. 4–11, New York. DOI: [10.1145/243199.243202](https://doi.org/10.1145/243199.243202). 144
- Y. Yang. 1999. An evaluation of statistical approaches to text categorization. *Journal of Information Retrieval*, 1:67–88. 315
- C. Zhai. 1997. Exploiting context to identify lexical atoms—a statistical view of linguistic context. In *Proc. of the International and Interdisciplinary Conference on Modelling and Using Context (CONTEXT-97)*, pages 119–129. Rio de Janeiro, Brazil. 273, 291
- C. Zhai. 2008. *Statistical Language Models for Information Retrieval*. Synthesis Lectures on Human Language Technologies. Morgan & Claypool Publishers. DOI: [10.2200/S00158ED1V01Y200811HLT001](https://doi.org/10.2200/S00158ED1V01Y200811HLT001). 55, 87, 128, 129
- C. Zhai and J. Lafferty. 2001. Model-based Feedback in the Language Modeling Approach to Information Retrieval. In *Proceedings of the Tenth International Conference on Information and Knowledge Management, CIKM '01*, ACM, pp. 403–410, New York. DOI: [10.1145/502585.502654](https://doi.org/10.1145/502585.502654). 143, 466, 473
- C. Zhai and J. Lafferty. 2004. A Study of Smoothing Methods for Language Models Applied to Information Retrieval. *ACM Trans. Inf. Syst.*, 22(2):179–214. 475
- C. Zhai, P. Jansen, E. Stoica, N. Grot, and D. A. Evans. 1998. Threshold Calibration in CLARIT Adaptive Filtering. In *Proc. of Seventh Text REtrieval Conference (TREC-7)*, pp. 149–156. 227

- C. Zhai, P. Jansen, and D. A. Evans. 2000. Exploration of a heuristic approach to threshold learning in adaptive filtering. In *SIGIR*, ACM, pp. 360–362. DOI: [10.1145/345508.345652](https://doi.org/10.1145/345508.345652). 235
- C. Zhai, A. Velivelli, and B. Yu. 2004. A cross-collection mixture model for comparative text mining. In *Proc. of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '04, ACM, pp. 743–748, New York. DOI: [10.1145/1014052.1014150](https://doi.org/10.1145/1014052.1014150). 423
- D. Zhang, C. Zhai, J. Han, A. Srivastava, and N. Oza. 2009. Topic modeling for OLAP on multidimensional text databases: topic cube and its applications. *Stat. Anal. Data Min.* 2, 5–6 (December 2009), 378–395. DOI: [10.1002/sam.v2.5/6](https://doi.org/10.1002/sam.v2.5/6). 440
- J. Zhu, A. Ahmed, and E. P. Xing. 2009. Medlda: Maximum margin supervised topic models for regression and classification. In *Proc. of the 26th Annual International Conference on Machine Learning*, ICML '09, ACM, pp. 1257–1264, New York. DOI: [10.1145/1553374.1553535](https://doi.org/10.1145/1553374.1553535). 386
- G. K. Zipf. 1949. *Human Behavior and the Principle of Least-Effort*. Cambridge, MA, Addison-Wesley, 162

Index

- Absolute discounting, 130
- Abstractive text summarization, 318, 321–324
- Access modes, 73–76
- Accuracy in search engine evaluation, 168
- Ad hoc information needs, 8–9
- Ad hoc retrieval, 75–76
- Add-1 smoothing, 130, 464
- Adjacency matrices, 207–208
- Advertising, opinion mining for, 393
- Agglomerative clustering, 277, 280–282, 290
- Aggregating
 - opinions, 393
 - scores, 234
- All-vs-all (AVA) method, 313
- Ambiguity
 - full structure parsing, 43
 - LARA, 406
 - NLP, 40–41, 44
 - one-vs-all method, 313
 - text retrieval vs. database retrieval, 80
 - topics, 335, 337
- Analyzers in MeTA toolkit, 61–64, 453
- `analyzers::filters` namespace, 64
- `analyzers::tokenizers` namespace, 64
- Anaphora resolution in natural language processing, 41
- Anchor text in web searches, 201
- Architecture
 - GFS, 194–195
 - MeTA toolkit, 60–61
 - unified systems, 452–453
- Art retrieval models, 111
- Aspect opinion analysis, 325–326
- Associations, word. *See* Word association mining
- Authority pages in web searches, 202, 207
- Automatic evaluation in text clustering, 294
- AVA (all-vs-all) method, 313
- Average-link document clustering, 282
- Average precision
 - ranked lists evaluation, 175, 177–180
 - search engine evaluation, 184
- Axiomatic thinking, 88
- Background models
 - mining topics from text, 345–351
 - mixture model estimation, 351–353
 - PLSA, 370–372
- Background words
 - mixture models, 141, 351–353
 - PLSA, 368–369, 372
- Bag-of-words
 - frequency analysis, 69
 - paradigmatic relations, 256
 - text information systems, 10
 - text representation, 88–90
 - vector space model, 93, 109
 - web searches, 215
- Bar-Hillel report, 42
- Baseline accuracy in text categorization, 314
- Bayes, Thomas, 25
- Bayes' rule
 - EM algorithm, 361–363, 373–374

- Bayes' rule (*continued*)
 - formula, 25–26
 - LDA, 383
- Bayesian inference
 - EM algorithm, 361–362
 - PLSA, 379, 382
- Bayesian parameter estimation
 - formula, 458
 - overfitting problem, 28–30
 - unigram language model, 341, 359
- Bayesian smoothing, 125
- Bayesian statistics
 - binomial estimation and beta distribution, 457–459
 - Dirichlet distribution, 461–463
 - LDA, 382
 - multinomial distribution, 460–461
 - multinomial parameters, 463–464
 - Naive Bayes algorithm, 309–312
 - pseudo counts, smoothing, and setting hyperparameters, 459–460
- Berkeley study, 3
- Bernoulli distribution, 26
- Beta distribution, 457–459
- Beta-gamma threshold learning, 227–228
- Bias, clustering, 276
- Big text data, 5–6
- Bigram language model
 - abstractive summarization, 323
 - Brown clustering, 290
- Bigrams
 - frequency analysis, 68
 - sentiment classification, 394–395
 - text categorization, 305
 - words tokenizers, 149
- Binary classification
 - content-based recommendation, 223
 - text categorization, 303
- Binary hidden variables in EM algorithm, 362–364, 366, 368, 467
- Binary logistic regression, 397
- Binomial distribution, 26–27
- Binomial estimation, 457–459
- Bit vector representation, 93–97
- Bitwise compression, 159–160
- Blind feedback, 133, 135
- Block compression, 161–162
- Block world project, 42
- BM25 model
 - description, 88
 - document clustering, 279
 - document length normalization, 108–109
 - link analysis, 201
 - Okapi, 89, 108
 - popularity, 90
 - probabilistic retrieval models, 111
- BM25-F model, 109
- BM25 score
 - paradigmatic relations, 258–261
 - syntagmatic relations, 270
 - web search ranking, 210
- BM25 TF transformation
 - description, 104–105
 - paradigmatic relations, 258–259
- BM25+ model, 88, 110
- Breadth-first crawler searches, 193
- Breakeven point precision, 189
- Brown clustering, 278, 288–291
- Browsing
 - multimode interactive access, 76–78
 - pull access mode, 73–75
 - support for, 445
 - text information systems, 9
 - web searches, 214
 - word associations, 252
- Business intelligence
 - opinion mining, 393
 - text data analysis, 243
- C++ language, 16, 58
- Caching
 - DBLRU, 164–165
 - LRU, 163–164
 - META toolkit, 60
 - search engine implementation, 148, 162–165
- Categories
 - categorical distributions, 460–461

- sentiment classification, 394, 396–397
- text information systems, 11–12
- Causal topic mining, 433–437
- Centroid vectors, 136–137
- Centroids in document clustering, 282–284
- CG (cumulative gain) in NDCG, 181–182
- character_tokenizer tokenizer, 61
- Citations, 202
- Classes
 - Brown clustering, 289
 - categories, 11–12
 - sentiment, 393–396
- Classification
 - machine learning, 34–36
 - NLP, 43–44
- Classifiers in text categorization, 302–303
- classify command, 57
- Cleaning HTML files, 218–219
- Clickthroughs
 - probabilistic retrieval model, 111–113
 - web searches, 201
- Clustering bias, 276
- Clusters and clustering
 - joint analysis, 416
 - sentiment classification, 395
 - text. *See* Text clustering
- Coherence in text clustering, 294–295
- Coin flips, binomial distribution for, 26–27
- Cold start problem, 230
- Collaborative filtering, 221, 229–233
- Collapsed Gibbs sampling, 383
- Collect function, 197
- Collection language model
 - KL-divergence, 474
 - smoothing methods, 121–126
- Common form of retrieval models, 88–90
- Common sense knowledge in NLP, 40
- Common words
 - background language model, 346–347, 350–351
 - feedback, 141, 143
 - filtering, 54
 - mixture models, 352–353, 355–356
 - unigram language model, 345–346
 - vector space retrieval models, 99, 109
- Compact clusters, 281
- Compare operator, 450, 452
- Complete data for EM algorithm, 467–468
- Complete-link document clustering, 281–282
- Component models
 - background language models, 345, 347–350
 - CPLSA, 421
 - description, 143
 - EM algorithm, 359
 - mixture models, 355–356, 358–359
 - PLSA, 370–373
- Compression
 - bitwise, 159–160
 - block, 161–162
 - overview, 158–159
 - search engines, 148
 - text representation, 48–49
- Compression ratio, 160–161
- Concepts in vector space model, 92
- Conceptual framework in text information systems, 10–13
- Conditional entropy
 - information theory, 33
 - syntagmatic relations, 261–264, 270
- Conditional probabilities
 - Bayes' rule, 25–26
 - overview, 23–25
- Configuration files, 57–58
- Confusion matrices, 314–315
- Constraints in PLSA, 373
- Content analysis modules, 10–11
- Content-based filtering, 221–229
- Content in opinion mining, 390–392
- Context
 - Brown clustering, 290
 - non-text data, 249
 - opinion mining, 390–392
 - paradigmatic relations, 253–258
 - social networks as, 428–433
 - syntagmatic relations, 261–262
 - text mining, 417–419

- Context (*continued*)
 - time series, 433–439
- Context variables in topic analysis, 330
- Contextual Probabilistic Latent Semantic Analysis (CPLSA), 419–428
- Continuous distributions
 - Bayesian parameter estimation, 28
 - description, 22
- Co-occurrences in mutual information, 267–268
- Corpus input formats in MeTA toolkit, 60–61
- corpusname.dat file, 60
- corpusname.dat.gz file, 60
- corpusname.dat.labels file, 60
- corpusname.dat.labels.gz file, 60
- Correlations
 - mutual information, 270
 - syntagmatic relations, 253–254
 - text-based forecasting, 248
 - time series context, 437
- Cosine similarity
 - document clustering, 279–280
 - extractive summarization, 321
 - text summarization, 325
 - vector measurement, 222, 232
- Coverage
 - CPLSA, 420–422, 425–426
 - LDA, 380–381
 - topic analysis, 332–333
- CPLSA (Contextual Probabilistic Latent Semantic Analysis), 419–428
- Cranfield evaluation methodology, 168–170
- Crawlers
 - domains, 218
 - dynamic content, 217
 - languages for, 216–217
 - web searches, 192–194
- Cross validation in text categorization, 314
- Cumulative gain (CG) in NDCG, 181–182
- Current technology, 5
- Data-driven social science research, opinion mining for, 393
- Data mining
 - joint analysis, 413–415
 - probabilistic retrieval model algorithms, 117
 - text data analysis, 245–246
- Data types in text analysis, 449–450
- Data-User-Service Triangle, 213–214
- Database retrieval, 80–82
- DBLRU (Double Barrel Least-Recently Used)
 - caches, 164–165
- DCG (discounted cumulative gain), 182–183
- Decision boundaries for linear classifiers, 311–312
- Decision modules in content-based filtering, 225
- Decision support, opinion mining for, 393
- Deep analysis in natural language processing, 43–45
- Delta bitwise compression, 160
- Dendrograms, 280–281
- Denial of service from crawlers, 193
- Dependency parsers, 323
- Dependent random variables, 25
- Design philosophy, MeTA, 58–59
- Development sets for text categorization, 314
- Dirichlet distribution, 461–463
- Dirichlet prior smoothing
 - KL-divergence, 475
 - probabilistic retrieval models, 125–127
- Disaster response, 243–244
- Discounted cumulative gain (DCG), 182–183
- Discourse analysis in NLP, 40
- Discrete distributions
 - Bayesian parameter estimation, 29
 - description, 22
- Discriminative classifiers, 302
- Distances in clusters, 281
- Distinguishing categories, 301–302
- Divergence-from-randomness models, 87, 111
- Divisive clustering, 277
- Document-at-a-time ranking, 155

- Document clustering, 277
 - agglomerative hierarchical, 280–282
 - K*-means, 282–284
 - overview, 279–280
- Document frequency
 - bag-of-words representation, 89
 - vector space model, 99–100
- Document IDs
 - compression, 158–159
 - inverted indexes, 152
 - tokenizers, 149
- Document language model, 118–123
- Document length
 - bag-of-words representation, 89
 - vector space model, 105–108
- Documents
 - filters, 155–156
 - ranking vs. selecting, 82–84
 - tokenizing, 148–150
 - vectors, 92–96
 - views in multimode interactive access, 77
- Domains, crawling, 218
- Dot products
 - document length normalization, 109
 - linear classifiers, 311
 - paradigmatic relations, 257–258
 - vector space model, 93–95, 98
- Double Barrel Least-Recently Used (DBLRU)
 - caches, 164–165
- Dynamic coefficient interpolation in
 - smoothing methods, 125
- Dynamically generated content and
 - crawlers, 217
- E step in EM algorithm, 362–368, 373–377, 465, 469
- E-discovery (electronic discovery), 326
- Edit features in text categorization, 306
- Effectiveness in search engine evaluation, 168
- Efficiency
 - database data retrieval, 81–82
 - search engine evaluation, 168
- Electronic discovery (E-discovery), 326
- Eliza project, 42, 44–45
- EM algorithm. *See* Expectation-maximization (EM) algorithm
- Email counts, 3
- Emotion analysis, 394
- Empirically defined problems, 82
- Enron email dataset, 326
- Entity-relation re-creation, 47
- Entropy
 - information theory, 31–33
 - KL-divergence, 139, 474
 - mutual information, 264–265
 - PMI, 288
 - skewed distributions, 158
 - syntagmatic relations, 261–264, 270
- Evaluation, search engine. *See* Search engine evaluation
- Events
 - CPLSA, 426–427
 - probability, 21–23
- Exhaustivity in sentiment classification, 396
- Expectation-maximization (EM) algorithm
 - CPLSA, 422
 - general procedure, 469–471
 - incomplete vs. complete data, 467–468
 - K*-means, 282–283
 - KL-divergence, 476
 - lower bound of likelihood, 468–469
 - MAP estimate, 378–379
 - mining topics from text, 359–368
 - mixture unigram language model, 466
 - MLE, 466–467
 - network supervised topic models, 431
 - overview, 465–466
 - PLSA, 373–377
- Expected overlap of words in paradigmatic
 - relations, 257–258
- Expected value in Beta distribution, 458
- Exploration-exploitation tradeoff in
 - content-based filtering, 227
- Extractive summarization, 318–321
- F measure
 - ranked lists evaluation, 179

- F measure (*continued*)
 - set retrieval evaluation, 172–173
- F*-test for time series context, 437
- F_1 score
 - text categorization, 314
 - text summarization, 324
- Fault tolerance in Google File System, 195
- Feature generation for tokenizers, 150
- Features for text categorization, 304–307
- Feedback
 - content-based filtering, 225
 - KL-divergence, 475–476
 - language models, 138–144
 - overview, 133–135
 - search engines, 147, 157–158
 - vector space model, 135–138
 - web searches, 201
- Feedback documents in unigram language model, 466
- Feelings. *See* Sentiment analysis
- `fetch_docs` function, 154
- `file_corpus` input format, 60
- Files in Google File System, 194–195
- Filter chains for tokenization, 61–64
- Filters
 - content-based, 221–229
 - documents, 155–156
 - recommender systems. *See* Recommender systems
 - text information systems, 11
 - unigram language models, 54
- Focused crawling, 193
- `forward_index` indexes, 60–61
- Forward indexes
 - description, 153
 - k*-nearest neighbors algorithm, 308
- Frame of reference encoding, 162
- Frequency and frequency counts
 - bag-of-words representation, 89–90
 - MapReduce, 197
 - META analyses, 68–70
 - term, 97–98
 - vector space model, 99–100
- Frequency transformation in paradigmatic relations, 258–259
- Full structure parsing, 43
- G*-means algorithm, 294
- Gain in search engine evaluation, 181–183
- Gamma bitwise compression, 160
- Gamma function, 457
- Gaussian distribution, 22, 404–405
- General EM algorithm, 431
- Generation-based text summarization, 318
- Generative classifiers, 309
- Generative models
 - background language model, 346–347, 349
 - CPLSA, 419, 421
 - description, 30, 36, 50
 - LARA, 403, 405–406
 - LDA, 381
 - log-likelihood functions, 343–344, 384
 - mining topics from text, 347
 - n*-gram models, 289
 - network supervised topic models, 428–430
 - PLSA, 370–371, 380
 - topics, 338–340
 - unigram language model, 341
- Geographical networks, 428
- Geometric mean average precision (gMAP), 179
- GFS (Google File System), 194–195
- Gibbs sampling, 383
- Google File System (GFS), 194–195
- Google PageRank, 202–206
- Grammar learning, 252
- Grammatical parse trees, 305–307
- Granger test, 434, 437
- Graph mining, 49
- `gz_corpus` input format, 60
- Hidden variables
 - EM algorithm, 362–364, 366, 368, 373–376, 465, 467
 - LARA, 403

- Hierarchical clustering, 280–282
- High-level syntactic features, 305–306
- Hill-climbing algorithm, EM, 360, 366–367, 465
- HITS algorithm, 206–208
- HTML files, cleaning, 218–219
- Hub pages in web searches, 202, 207–208
- Humans
 - joint analysis, 413–415
 - NLP, 48
 - opinion mining. *See* Opinion mining
 - as subjective sensors, 244–246
 - unified systems, 445–448
- Hyperparameters
 - Beta distribution, 458–460
 - Dirichlet distribution, 461, 463
- ICU (International Components for Unicode), 61
- Icu_filter filter, 61
- Icu_tokenizer tokenizer, 61
- IDF (inverse document frequency)
 - Dirichlet prior smoothing, 126
 - paradigmatic relations, 258–260
 - query likelihood retrieval model, 122
 - vector space model, 99–101
- Illinois NLP Curator toolkit, 64
- Impact
 - CPLSA, 426–427
 - time series context, 437
- Implicit feedback, 134–135
- Incomplete data in EM algorithm, 467–468
- Incremental crawling, 193
- Independent random variables, 25
- Index sharding, 156–157
- Indexes
 - compressed, 158–162
 - forward, 153, 308
 - k*-nearest neighbors algorithm, 308
 - MapReduce, 198–199
 - META toolkit, 60–61, 453–455
 - search engine implementation, 150–153
 - search engines, 147, 150–153
 - text categorization, 314
 - web searches, 194–200
- Indirect citations in web searches, 202
- Indirect opinions, 391–392
- Indri/Lemur search engine toolkit, 64
- Inferences
 - NLP, 41
 - probabilistic, 88
 - real world properties, 248
- Inferred opinions, 391–392
- Information access in text information systems, 7
- Information extraction
 - NLP, 43
 - text information systems, 9, 12
- Information retrieval (IR) systems, 6
 - evaluation metrics, 324–325
 - implementation. *See* Search engine implementation
 - text data access, 79
- Information theory, 31–34
- Initial values in EM algorithm, 466
- Initialization modules in content-based filtering, 224–225
- Inlink counts in PageRank, 203
- Instance-based classifiers, 302
- Instructor reader category, 16–17
- Integer compression, 158–162
- Integration of information access in web searches, 213
- Integrity in text data access, 81
- Interactive access, multimode, 76–78
- Interactive task support in web searches, 216
- International Components for Unicode (ICU), 61
- Interpolation for smoothing methods, 125–126
- Interpret operator, 450–452
- Intersection operator, 449–450
- Intrusion detection, 271–273
- Inverse document frequency (IDF)
 - Dirichlet prior smoothing, 126
 - paradigmatic relations, 258–260
 - query likelihood retrieval model, 122

- Inverse document frequency (IDF)
 - (*continued*)
 - vector space model, 99–101
- Inverse user frequency (IUF), 232
- inverted_index indexes, 60
- Inverted index chunks, 156–157
- Inverted indexes
 - compression, 158
 - k -nearest neighbors algorithm, 308
 - MapReduce, 198–199
 - search engines, 150–153
- IR (information retrieval) systems, 6
 - evaluation metrics, 324–325
 - implementation. *See* Search engine implementation
 - text data access, 79
- Iterative algorithms for PageRank, 205–206
- Iterative Causal Topic Modeling, 434–435
- IUF (inverse user frequency), 232
- Jaccard similarity, 280
- Jelinek-Mercer smoothing, 123–126
- Joint analysis of text and structured data, 413
 - contextual text mining, 417–419
 - CPLSA, 419–428
 - introduction, 413–415
 - social networks as context, 428–433
 - time series context, 433–439
- Joint distributions for mutual information, 266–268
- Joint probabilities, 23–25
- K -means document clustering, 282–284
- K -nearest neighbors (k -NN) algorithm, 307–309
- Kernel trick for linear classifiers, 312
- Key-value pairs in MapReduce, 195–198
- KL-divergence
 - Dirichlet prior smoothing, 475
 - EM algorithm, 468
 - feedback, 139–140
 - mutual information, 266
 - query model, 475–476
 - retrieval, 473–474
- Knowledge acquisition in text information systems, 8–9
- Knowledge discovery in text summarization, 326
- Knowledge Graph, 215
- Knowledge provenance in unified systems, 447
- Known item searches in ranked lists
 - evaluation, 179
- Kolmogorov axioms, 22–23
- Kullback-Leibler divergence retrieval model. *See* KL-divergence
- Lagrange Multiplier approach
 - EM algorithm, 467, 470
 - unigram language model, 344
- Language models
 - feedback in, 138–144
 - in probabilistic retrieval model, 87, 111, 117
- Latent Aspect Rating Analysis (LARA), 400–409
- Latent Dirichlet Allocation (LDA), 377–383
- Latent Rating Regression, 402–405
- Lazy learners in text categorization, 302
- Learners
 - search engines, 147
 - text categorization, 302
- Learning modules in content-based filtering, 224–225
- Least-Recently Used (LRU) caches, 163–164
- length_filter filter, 61
- Length normalization
 - document length, 105–108
 - query likelihood retrieval model, 122
- Lexical analysis in NLP, 39–40
- Lexicons for inverted indexes, 150–152
- LIBLINEAR algorithm, 58
- libsvm_analyzer analyzer, 62
- libsvm_corpus file, 61
- LIBSVM package, 58, 64
- Lifelong learning in web searches, 213

- Likelihood and likelihood function
 - background language model, 349–351
 - EM algorithm, 362–363, 367–368, 376, 465–469
 - LARA, 405
 - LDA, 378, 381–382
 - marginal, 28
 - mixture model behavior, 354–357
 - MLE, 27
 - network supervised topic models, 428–431
 - PLSA, 372–374
 - unigram language model, 342–344
- line_corpus input format, 60
- Linear classifiers in text categorization, 311–313
- Linear interpolation in Jelinek-Mercer smoothing, 124
- Linearly separable data points in linear classifiers, 312
- Link analysis
 - HITS, 206–208
 - overview, 200–202
 - PageRank, 202–206
- list_filter filter, 62
- Local maxima, 360, 363, 367–368, 465
- Log-likelihood function
 - EM algorithm, 365–366, 466–467
 - feedback, 142–143
 - unigram language model, 343–344
- Logarithm transformation, 103–104
- Logarithms in probabilistic retrieval model, 118, 122
- Logic-based approach in NLP, 42
- Logical predicates in NLP, 49–50
- Logistic regression in sentiment classification, 396–400
- Long-range jumps in multimode interactive access, 77
- Long-term needs in push access mode, 75
- Low-level lexical features in text categorization, 305
- Lower bound of likelihood in EM algorithm, 468–469
- LRU (Least-Recently Used) caches, 163–164
- Lucene search engine toolkit, 64
- M step
 - EM algorithm, 361–368, 373–377, 465, 469–470
 - MAP estimate, 379
 - network supervised topic models, 431
- Machine-generated data, 6
- Machine learning
 - overview, 34–36
 - sentiment classification methods, 396
 - statistical, 10
 - text categorization, 301
 - web search algorithms, 201
 - web search ranking, 208–212
- Machine translation, 42, 44–45
- Magazine output, 3
- Manual evaluation for text clustering, 294
- map function, 195–198
- MAP (Maximum a Posteriori) estimate
 - Bayesian parameter estimation, 29
 - LARA, 404–405
 - PLSA, 378–379
 - word association mining, 271–273
- MAP (mean average precision), 178–180
- Map Reduce paradigm, 157
- MapReduce framework, 194–200
- Maps in multimode interactive access, 76–77
- Marginal probabilities
 - Bayesian parameter estimation, 29
 - mutual information, 267
- Market research, opinion mining for, 393
- Massung, Sean, biography, 490
- Matrices
 - adjacency, 207–208
 - PageRank, 204–208
 - text categorization, 314–315
 - transition, 204
- Matrix multiplication in PageRank, 205
- Maximal marginal relevance (MMR)
 - reranking
 - extractive summarization, 320–321

- Maximal marginal relevance (MMR)
 - reranking (*continued*)
 - topic analysis, 333
- Maximization algorithm for document clustering, 282
- Maximum a Posteriori (MAP) estimate
 - Bayesian parameter estimation, 29
 - LARA, 404–405
 - PLSA, 378–379
 - word association mining, 271–273
- Maximum likelihood estimation (MLE)
 - background language model, 346, 350
 - Brown clustering, 289
 - Dirichlet prior smoothing, 125–126
 - EM algorithm, 359–368, 466–467
 - feedback, 141–143
 - generative models, 339
 - Jelinek-Mercer smoothing, 124
 - KL-divergence, 475–476
 - LARA, 404
 - LDA, 382
 - mixture model behavior, 354–359
 - mixture model estimation, 352–353
 - multinomial distribution, 463
 - mutual information, 268–269
 - overview, 27–28
 - PLSA, 372–373, 378
 - query likelihood retrieval model, 118–119
 - term clustering, 286
 - unigram language models, 52–53, 341–345
 - web search ranking, 210
- Mean average precision (MAP), 178–180
- Mean reciprocal rank (MRR), 180
- Measurements in search engine evaluation, 168
- Memory-based approach in collaborative filtering, 230
- META toolkit
 - architecture, 60–61
 - classification algorithms, 307
 - design philosophy, 58–59
 - exercises, 65–70
 - overview, 57–58
 - related toolkits, 64–65
 - setting up, 59–60
 - text categorization, 314–315
 - tokenization, 61–64
 - as unified system, 453–455
- Metadata
 - classification algorithms, 307
 - contextual text mining, 417
 - networks from, 428
 - text data analysis, 249
 - topic analysis, 330
- Mining
 - contextual, 417–419
 - demand for, 4–5
 - graph, 49
 - joint analysis, 413–419
 - opinion. *See* Opinion mining; Sentiment analysis
 - probabilistic retrieval model, 117
 - tasks, 246–250
 - toolkits, 64
 - topic analysis, 330–331
 - word association. *See* Word association mining
- Mining topics from text, 340
 - background language model, 345–351
 - expectation-maximization, 359–368
 - joint analysis, 416
 - mixture model behavior, 353–359
 - mixture model estimation, 351–353
 - unigram language model, 341–345
- Mixture models
 - behavior, 353–359
 - EM algorithm, 466
 - estimation, 351–353
 - feedback, 140–142, 157
 - mining topics from text, 346–351
- MLE. *See* Maximum likelihood estimation (MLE)
- MMR (maximal marginal relevance)
 - reranking
 - extractive summarization, 320–321
 - topic analysis, 333
- Model-based clustering algorithms, 276–277
- Model files for META toolkit, 59

- Modification in NLP, 41
- Modules in content-based filtering, 224–226
- MRR (mean reciprocal rank), 180
- Multiclass classification
 - linear classifiers, 313
 - text categorization, 303
- Multi-level judgments in search engine
 - evaluation, 180–183
- Multimode interactive access, 76–78
- Multinomial distributions
 - Bayesian estimate, 463–464
 - generalized, 460–461
 - LDA, 380
- Multinomial parameters in Bayesian
 - estimate, 463–464
- Multiple-level sentiment analysis, 397–398
- Multiple occurrences in vector space model, 103–104
- Multiple queries in ranked lists evaluation, 178–180
- Multivariate Gaussian distribution, 404–405
- Mutual information
 - information theory, 33–34
 - syntagmatic relations, 264–271
 - text clustering, 278
- n*-fold cross validation, 314
- n*-gram language models
 - abstractive summarization, 322–323
 - frequency analysis, 68–69
 - sentiment classification, 394–395
 - term clustering, 288–291
 - vector space model, 109
- Naive Bayes algorithm, 309–312
- Named entity recognition, 323
- Natural language, mining knowledge about, 247
- Natural language generation in text
 - summarization, 323–324
- Natural language processing (NLP)
 - history and state of the art, 42–43
 - pipeline, 306–307
 - sentiment classification, 395
 - statistical language models, 50–54
 - tasks, 39–41
 - text information systems, 43–45
 - text representation, 46–50
- Navigating maps in multimode interactive
 - access, 77
- Navigational queries, 200
- NDCG (normalized discounted cumulative gain), 181–183
- NDCG@*k* score, 189
- Nearest-centroid classifiers, 309
- Negative feedback documents, 136–138
- Negative feelings, 390–394
- NetPLSA model, 430–433
- Network supervised topic models, 428–433
- Neural language model, 291–294
- News summaries, 317
- Newspaper output, 3
- ngram_pos_analyzer* analyzer, 62
- ngram_word_analyzer* analyzer, 62
- NLP. *See* Natural language processing (NLP)
- NLTK toolkit, 64
- no_evict_cache* caches, 60
- Nodes in word associations, 252
- Non-text data
 - context, 249
 - predictive analysis, 249
 - vs. text, 244–246
- Normalization
 - document length, 105–108
 - PageRank, 206
 - query likelihood retrieval model, 122
 - term clustering, 286
 - topic analysis, 333
- Normalized discounted cumulative gain (NDCG), 181–183
- Normalized ratings in collaborative
 - filtering, 230–231
- Normalized similarity algorithm, 279
- Objective statements vs. subjective, 389–390
- Observed world, mining knowledge about, 247–248
- Observers, mining knowledge about, 248
- Office documents, 3
- Okapi BM25 model, 89, 108
- One-vs-all (OVA) method, 313

- Operators in text analysis systems, 448–452
- Opinion analysis in text summarization, 325–326
- Opinion holders, 390–392
- Opinion mining
 - evaluation, 409–410
 - LARA, 400–409
 - overview, 389–392
 - sentiment classification. *See* Sentiment analysis
- Opinion summarization, 318
- Optimization in web searches, 191
- Ordinal regression, 394, 396–400
- Organization in text information systems, 8
- OVA (one-vs-all) method, 313
- Over-constrained queries, 84
- Overfitting problem
 - Bayesian parameter estimation, 28, 30
 - sentiment classification, 395
 - vector space model, 138
- Overlap of words in paradigmatic relations, 257–258
- p*-values in search engine evaluation, 185–186
- PageRank technique, 202–206
- Paradigmatic relations
 - Brown clustering, 290
 - discovering, 252–260
 - overview, 251–252
- Parallel crawling, 193
- Parallel indexing and searching, 192
- Parameters
 - background language model, 350–351
 - Bayesian parameter estimation, 28–30, 341, 359, 458, 463–464
 - Beta distribution, 458–460
 - Dirichlet distribution, 461–463
 - EM algorithm, 363, 465
 - feedback, 142–144
 - LARA, 404–405
 - LDA, 380–381
 - mixture model estimation, 352
 - MLE. *See* Maximum likelihood estimation (MLE)
 - network supervised topic models, 429
 - PLSA, 372–373, 379–380
 - probabilistic models, 30–31
 - ranking, 209–211
 - statistical language models, 51–52
 - topic analysis, 338–339
 - unigram language models, 52
- Parsing
 - META toolkit, 67–68
 - NLP, 43
 - web content, 216
- Part-of-speech (POS) tags
 - META toolkit, 67
 - NLP, 47
 - sentiment classification, 395
- Partitioning
 - Brown clustering, 289
 - extractive summarization, 319–320
 - text data, 417–419
- Patterns
 - contextual text mining, 417–419
 - CPLSA, 425–426
 - joint analysis, 417
 - NLP, 45
 - sentiment classification, 395
- Pdf (probability density function)
 - Beta distribution, 457
 - Dirichlet distribution, 461
 - multinomial distribution, 461
- Pearson correlation
 - collaborative filtering, 222, 231–232
 - time series context, 437
- Perceptron classifiers, 312–313
- Personalization in web searches, 212, 215
- Personalized PageRank, 206
- Perspective in text data analysis, 246–247
- Pivoted length normalization, 89, 107–108
- PL2 model, 90
- PLSA (probabilistic latent semantic analysis)
 - CPLSA, 419–428
 - extension, 377–383
 - overview, 368–377
- Pointwise Mutual Information (PMI), 278, 287–288

- Polarity analysis in sentiment classification, 394
- Policy design, opinion mining for, 393
- Pooling in search engine evaluation, 186–187
- Porter2 English Stemmer, 66–67
- porter2_stemmer filter, 62
- POS (part-of-speech) tags
 - META toolkit, 67
 - NLP, 47
 - sentiment classification, 395
- Positive feelings, 390–394
- Posterior distribution, 28
- Posterior probability in Bayesian parameter estimation, 29
- Postings files for inverted indexes, 150–152
- Power iteration for PageRank, 205
- Practitioners reader category, 17
- Pragmatic analysis in NLP, 39–40
- Precision
 - search engine evaluation, 184
 - set retrieval evaluation, 170–178
- Precision-recall curves in ranked lists
 - evaluation, 174–176
- Predictive analysis for non-text data, 249
- Predictors features in joint analysis, 413–416
- Presupposition in NLP, 41
- Prior probability in Bayesian parameter estimation, 29
- Probabilistic inference, 88
- Probabilistic latent semantic analysis (PLSA)
 - CPLSA, 419–428
 - extension, 377–383
 - overview, 368–377
- Probabilistic retrieval models
 - description, 87–88
 - overview, 110–112
 - query likelihood retrieval model, 114–118
- Probability and statistics
 - abstractive summarization, 322
 - background language model, 346–349
 - basics, 21–23
 - Bayes' rule, 25–26
 - Bayesian parameter estimation, 28–30
 - binomial distribution, 26–27
 - EM algorithm, 362–366
 - joint and conditional probabilities, 23–25
 - KL-divergence, 474
 - LARA, 403
 - maximum likelihood parameter estimation, 27–28
 - mixture model behavior, 354–358
 - mutual information, 266–270
 - Naive Bayes algorithm, 310
 - PageRank, 202–206
 - paradigmatic relations, 257–258
 - PLSA, 368–377, 380
 - probabilistic models and applications, 30–31
 - syntagmatic relations, 262–263
 - term clustering, 286–289
 - topics, 336–339
 - unigram language model, 342–344
 - web search ranking, 209–211
- Probability density function (pdf)
 - Beta distribution, 457
 - Dirichlet distribution, 461
 - multinomial distribution, 461
- Probability distributions
 - overview, 21–23
 - statistical language models, 50–54
- Probability ranking principle, 84
- Probability space, 21–23
- Producer-initiated recommendations, 75
- Product reviews in opinion mining, 391–392
- profile command, 65–66
- Properties
 - inferring knowledge about, 248
 - text categorization for, 300
- Proximity heuristics for inverted indexes, 151
- Pseudo counts
 - Bayesian statistics, 459–460
 - LDA, 381
 - multinomial distribution, 463
 - PLSA, 379, 381
 - smoothing techniques, 128, 286
- Pseudo data in LDA, 378

- Pseudo feedback, 133, 135, 142, 157–158
- Pseudo-segments for mutual information, 269–270
- Pull access mode, 8–9, 73–76
- Push access mode, 8–9, 73–76
- Python language
 - cleaning HTML files, 218
 - crawlers, 217
- Q-function, 465, 469–471
- Queries
 - multimode interactive access, 77
 - navigational, 200
 - text information systems, 9
 - text retrieval vs. database retrieval, 80
- Query expansion
 - vector space model, 135
 - word associations, 252
- Query likelihood retrieval model, 90, 113
 - document language model, 118–123
 - feedback, 139
 - KL-divergence, 475–476
 - overview, 114–118
 - smoothing methods, 123–128
- Query vectors, 92–98, 135–137
- Random access decoding in compression, 158
- Random numbers in abstractive summarization, 322
- Random observations in search engine evaluation, 186
- Random surfers in PageRank, 202–204
- Random variables
 - Bayesian parameter estimation, 28
 - dependent, 25
 - entropy of, 158, 262–263, 270
 - information theory, 31–34
 - PMI, 287
 - probabilistic retrieval models, 87, 111, 113
 - probability distributions, 22
- Ranked lists evaluation
 - multiple queries, 178–180
 - overview, 174–178
- Rankers for search engines, 147
- Ranking
 - extractive summarization, 320
 - probabilistic retrieval model. *See* Probabilistic retrieval models
 - vs. selection, 82–84
 - text analysis operator, 450–451
 - text data access, 78
 - vector space model. *See* Vector space (VS) retrieval models
 - web searches, 201, 208–212
- Ratings
 - collaborative filtering, 230–231
 - LARA, 400–409
 - sentiment classification, 396–399
- Real world properties, inferring knowledge about, 248
- Realization in abstractive summarization, 324
- Recall in set retrieval evaluation, 170–178
- Reciprocal ranks, 179–180
- Recommendations in text information systems, 11
- Recommender systems
 - collaborative filtering, 229–233
 - content-based recommendation, 222–229
 - evaluating, 233–235
 - overview, 221–222
- reduce function, 198
- Redundancy
 - MMR reranking, 333
 - text summarization, 320–321, 324
 - vector space retrieval models, 92
- Regression
 - LARA, 402–405
 - machine learning, 34–35
 - sentiment classification, 394, 396–400
 - text categorization, 303–304
 - web search ranking, 209–211
- Regularizers in network supervised topic models, 429–431
- Relevance and relevance judgments
 - Cranfield evaluation methodology, 168–169

- description, 133
- document ranking, 83
- document selection, 83
- extractive summarization, 321
- probabilistic retrieval models, 110–112
- search engine evaluation, 181–184, 186–187
- set retrieval evaluation, 171–172
- text data access, 79
- vector space model, 92
- web search ranking, 209–211
- Relevant text data, 5–6
- Relevant word counts in EM algorithm, 364–365, 376–377
- Repeated crawling, 193
- Representative documents in search engine evaluation, 183
- reset command, 57–59
- Retrieval models
 - common form, 88–90
 - overview, 87–88
 - probabilistic. *See* Probabilistic retrieval models
 - vector space. *See* Vector space (VS) retrieval models
- Reviews
 - LARA, 400–409
 - opinion mining, 391–392
 - sentiment classification, 394
 - text summarization, 318
- RMSE (root-mean squared error), 233
- robots.txt file, 193
- Rocchio feedback
 - forward indexes, 157
 - vector space model, 135–138
- Root-mean squared error (RMSE), 233
- Ruby language
 - cleaning html files, 218–219
 - crawlers, 217
- Rule-based text categorization, 301
- Scalability in web searches, 191–192
- Scanning inverted indexes, 152
- Scientific research, text data analysis for, 243
- Scikit Learn toolkit, 64
- score_term function, 154
- Scorers
 - document-at-a-time ranking, 155
 - filtering documents, 155–156
 - index sharding, 156–157
 - search engines, 147, 153–157
 - term-at-a-time ranking, 154–155
- Scoring functions
 - KL-divergence, 474
 - topic analysis, 332
- SDI (selective dissemination of information), 75
- Search engine evaluation
 - Cranfield evaluation methodology, 168–170
 - measurements, 168
 - multi-level judgments, 180–183
 - practical issues, 183–186
 - purpose, 167–168
 - ranked lists, 174–180
 - set retrieval, 170–173
- Search engine implementation
 - caching, 162–165
 - compression, 158–162
 - feedback implementation, 157–158
 - indexes, 150–153
 - overview, 147–148
 - scorers, 153–157
 - tokenizers, 148–150
- Search engine queries
 - pull access mode, 74–75
 - text data access, 78–79
- Search engine toolkits, 64
- Searches
 - text information systems, 11
 - web. *See* Web searches
- Segmentation in LARA, 405
- Select operator, 449–451, 455
- Selection
 - vs. ranking, 82–84
 - text data access, 78
- Selection-based text summarization, 318
- Selective dissemination of information (SDI), 75

- Semantic analysis in NLP, 39–40, 43, 47
- Semantically related terms in clustering, 187, 285–287
- Sensors
 - humans as, 244–246
 - joint analysis, 413–415
 - opinion mining. *See* Opinion mining
- Sentence vectors in extractive summarization, 319
- Sentiment analysis, 389
 - classification, 393–396
 - evaluation, 409–410
 - NLP, 43
 - ordinal regression, 396–400
 - text categorization, 304
- Separation in text clustering, 294–295
- Sequences of words in NLP, 46–47
- Set retrieval evaluation
 - description, 170
 - F measure, 172–173
 - precision and recall, 170–173
- Shadow analysis in NLP, 48
- Shallow analysis in NLP, 43–45
- Short-range walks in multimode interactive access, 77
- Short-term needs in pull access mode, 75
- Sign tests in search engine evaluation, 185
- Signed-rank tests in search engine evaluation, 185
- Significance tests in search engine evaluation, 183–186
- Similarity algorithm for clustering, 276
- Similarity in clustering
- Similarity functions and measures
 - extractive summarization, 319, 321
 - paradigmatic relations, 256–259
 - vector space model, 92, 109
 - description, 277
 - document clustering, 279–281
 - term clustering, 285
- Single-link document clustering, 281–282
- Skip-gram neural language model, 292–293
- sLDA (supervised LDA), 387
- Smoothing techniques
 - Add-1, 130, 464
 - Bayesian statistics, 459–460
 - KL-divergence, 474–475
 - maximum likelihood estimation, 119–128
 - multinomial distribution, 463–464
 - Naive Bayes algorithm, 310
 - unigram language models, 53
- Social media in text data analysis, 243
- Social networks as context, 428–433
- Social science research, opinion mining for, 393
- Soft rules in text categorization, 301
- Spam in web searches, 191–192
- Sparse Beta, 459
- Sparse data in Naive Bayes algorithm, 309–311
- Sparse priors in Dirichlet distribution, 461–462
- Spatiotemporal patterns in CPLSA, 425–426
- Specificity in sentiment classification, 396
- Speech acts in NLP, 47–48
- Speech recognition
 - applications, 42
 - statistical language models, 51
- Spiders for web searches, 192–194
- Split counts in EM algorithm, 374–375
- Split operator for text analysis, 449–452, 455
- Stanford NLP toolkit, 64
- State-of-the-art support vector machines (SVM) classifiers, 311–312
- Statistical language models
 - NLP, 45
 - overview, 50–54
- Statistical machine learning
 - NLP, 42–43
 - text information systems, 10
- Statistical significance tests in search engine evaluation, 183–186
- Statistics. *See* Probability and statistics
- Stemmed words in vector space model, 109
- Stemming process in META toolkit, 66–67
- Sticky phrases in Brown clustering, 291
- Stop word removal
 - feedback, 141
 - frequency analysis, 69

- META toolkit, 62, 66
 - mixture models, 352
 - vector space model, 99, 109
- Story understanding, 42
- Structured data
 - databases, 80
 - joint analysis with text. *See* Joint analysis of text and structured data
- Student reader category, 16
- Stylistic analysis in NLP, 49
- Subjective sensors
 - humans as, 244–246
 - opinion mining. *See* Opinion mining
- Subjective statements vs. objective, 389–390
- Sublinear transformation
 - term frequency, 258–259
 - vector space model, 103–104
- Summarization. *See* Text summarization
- Supervised LDA (sLDA), 387
- Supervised machine learning, 34
- SVM (state-of-the-art support vector machines) classifiers, 311–312
- Symbolic approach in NLP, 42
- Symmetric Beta, 459
- Symmetric probabilities in information theory, 32
- Symmetry in document clustering, 279–280
- Synonyms
 - vector space model, 92
 - word association, 252
- Syntactic ambiguity in NLP, 41
- Syntactic analysis in NLP, 39–40, 47
- Syntactic structures in NLP, 49
- SyntacticDiff method, 306
- Syntagmatic relations, 251–252
 - Brown clustering, 290–291
 - discovering, 253–254, 260–264
 - mutual information, 264–271
- System architecture in unified systems, 452–453
- Tags, POS
 - META toolkit, 67
 - NLP, 47
 - sentiment classification, 395
- Targets in opinion mining, 390–392
- Temporal trends in CPLSA, 424–425
- Term-at-a-time ranking, 154–155
- Term clustering, 278
 - n*-gram class language models, 288–291
 - neural language model, 291–294
 - overview, 284–285
 - Pointwise Mutual Information, 287–288
 - semantically related terms, 285–287
- Term frequency (TF)
 - bag-of-words representation, 89
 - vector space model, 97–98
- Term IDs
 - inverted indexes, 151–152
 - tokenizers, 149–150
- Term vectors, 92
- Terms, topics as, 332–335
- Terrier search engine toolkit, 64
- Test collections in Cranfield evaluation
 - methodology, 168–169
- Testing data
 - machine learning, 35
 - text categorization, 303
- Text
 - joint analysis with structured data. *See* Joint analysis of text and structured data
 - mining. *See* Mining; Mining topics from text
 - usefulness, 3–4
- Text annotation. *See* Text categorization
- Text-based prediction, 300
- Text categorization
 - classification algorithms overview, 307
 - evaluation, 313–315
 - features, 304–307
 - introduction, 299–301
 - k*-nearest neighbors algorithm, 307–309
 - linear classifiers, 311–313
 - machine learning, 35
 - methods, 300–302
 - Naive Bayes, 309–311
 - problem, 302–304
- Text clustering, 12
 - document, 279–284

- Text clustering (*continued*)
 - evaluation, 294–296
 - overview, 275–276
 - techniques, 277–279
 - term, 284–294
- Text data access, 73
 - access modes, 73–76
 - document selection vs. document ranking, 82–84
 - multimode interactive, 76–78
 - text retrieval vs. database retrieval, 80–82
 - text retrieval overview, 78–80
- Text data analysis overview, 241–242
 - applications, 242–244
 - humans as subjective sensors, 244–246
 - operators, 448–452
 - text information systems, 8
 - text mining tasks, 246–250
- Text data understanding. *See* Natural language processing (NLP)
- Text information systems (TISs)
 - conceptual framework, 10–13
 - functions, 7–10
 - NLP, 43–45
- Text management and analysis in unified systems. *See* Unified systems
- Text organization in text information systems, 8
- Text representation in NLP, 46–50
- Text retrieval (TR)
 - vs. database retrieval, 80–82
 - demand for, 4–5
 - overview, 78–80
- Text summarization, 12
 - abstractive, 321–324
 - applications, 325–326
 - evaluation, 324–325
 - extractive, 319–321
 - overview, 317–318
 - techniques, 318
- TextObject data type operators, 449, 454
- TextObjectSequence data type operators, 449, 454
- TF (term frequency)
 - bag-of-words representation, 89
 - vector space model, 97–98
- TF-IDF weighting
 - Dirichlet prior smoothing, 128
 - probabilistic retrieval model, 122–123
 - topic analysis, 333
 - vector space model, 100–103
- TF transformation, 102–105
- TF weighting, 125–126
- Themes in CPLSA, 420–422
- Therapist application, 44–45
- Thesaurus discovery in NLP, 49
- Threshold settings in content-based filtering, 222, 224–227
- Tight clusters, 281
- Time series context in topic analysis, 433–439
- TISs (text information systems)
 - conceptual framework, 10–13
 - functions, 7–10
 - NLP, 43–45
- Tokenization
 - META toolkit, 61–64, 453
 - search engines, 147–150
- Topic analysis
 - evaluation, 383–384
 - LDA, 377–383
 - mining topics from text. *See* Mining topics from text
 - model summary, 384–385
 - overview, 329–331
 - PLSA, 368–377
 - social networks as context, 428–433
 - text information systems, 12
 - time series context, 433–439
 - topics as terms, 332–335
 - topics as word distributions, 335–340
- Topic coherence in time series context, 436
- Topic coverage
 - CPLSA, 420–422, 425–426
 - LDA, 380–381
- Topic maps in multimode interactive access, 76–77
- TopicExtraction operator, 450
- TR (text retrieval)
 - vs. database retrieval, 80–82

- demand for, 4–5
 - overview, 78–80
- Training and training data
 - classification algorithms, 307–309
 - collaborative filtering, 229–230
 - content-based recommendation, 227–228
 - linear classifiers, 311–313
 - machine learning, 34–36
 - Naive Bayes, 309–310
 - NLP, 42–43, 45
 - ordinal regression, 398–399
 - text categorization, 299–303, 311–314
 - web search ranking, 209–210, 212
- Transformations
 - frequency, 258–259
 - vector space model, 103–104
- Transition matrices in PageRank, 204
- Translation, machine, 42, 44–45
- TREC filtering tasks, 228
- tree_analyzer analyzer, 62
- Trends in web searches, 215–216
- Trigrams in frequency analysis, 69
- Twitter searches, 83
- Two-component mixture model, 356
- Unary bitwise compression, 159–160
- Under-constrained queries, 84
- Unified systems
 - MeTA as, 453–455
 - overview, 445–448
 - system architecture, 452–453
 - text analysis operators, 448–452
- Uniform priors in Dirichlet distribution, 461
- Unigram language models, 51–54
 - EM algorithm, 466
 - LDA, 381
 - mining topics from text, 341–345
 - PLSA, 370
- Unigrams
 - abstractive summarization, 321–323
 - frequency analysis, 68
 - sentiment classification, 394
 - words tokenizers, 149
- Unimodel Beta, 459
- Union operator, 449–450
- University of California Berkeley study, 3
- Unseen words
 - document language model, 119–120, 122–123
 - KL-divergence, 474
 - Naive Bayes algorithm, 310–311
 - smoothing, 124, 285–287
 - statistical language models, 52
- Unstructured text access, 80
- Unsupervised clustering algorithms, 275, 278
- Unsupervised machine learning, 34, 36
- URLs and crawlers, 193
- Usability in search engine evaluation, 168
- Utility
 - content-based filtering, 224–228
 - text clustering, 294
- Valence scoring, 411
- Variable byte encoding, 161
- Variables
 - context, 330
 - contextual text mining, 419
 - CPLSA, 422
 - EM algorithm, 362–364, 366, 368, 373–376, 465, 467
 - LARA, 403
 - random. *See* Random variables
- vByte encoding, 161
- Vector space (VS) retrieval models, 87
 - bit vector representation, 94–97
 - content-based filtering, 225–226
 - document length normalization, 105–108
 - feedback, 135–138
 - improved instantiation, 97–102
 - improvement possibilities, 108–110
 - instantiation, 93–95
 - overview, 90–92
 - paradigmatic relations, 256–258
 - summary, 110
 - TF transformation, 102–105
- Vectors
 - collaborative filtering, 222

- Vectors (*continued*)
 - neural language model, 292
- Versions, META toolkit, 59
- Vertical search engines, 212
- Video data mining, 245
- Views
 - CPLSA, 420–422
 - multimode interactive access, 77
- Visualization in text information systems, 12–13
- VS retrieval models. *See* Vector space (VS) retrieval models
- Web searches
 - crawlers, 192–194
 - future of, 212–216
 - indexing, 194–200
 - link analysis, 200–208
 - overview, 191–192
 - ranking, 208–212
- Weighted k -nearest neighbors algorithm, 309
- WeightedTextObjectSequence data type, 449
- WeightedTextObjectSet data type, 449
- Weights
 - collaborative filtering, 231
 - Dirichlet prior smoothing, 127–128
 - document clustering, 279–280
 - LARA, 401–409
 - linear classifiers, 313
 - mutual information, 269–270
 - NetPLSA model, 430
 - network supervised topic models, 431
 - paradigmatic relations, 258–261
 - query likelihood retrieval model, 121–123
 - text categorization rules, 301
 - topics, 333, 335–336
 - vector space model, 92, 99–103
- Weka toolkit, 64
- whitespace_tokenizer command, 149
- Whitespace tokenizers, 149
- Wilcoxon signed-rank test, 185
- Word association mining
 - evaluation, 271–273
 - general idea of, 252–254
 - overview, 251–252
 - paradigmatic relations discovery, 254–260
 - syntagmatic relations discovery, 260–271
- Word counts
 - EM algorithm, 364–365, 376–377
 - MapReduce, 195–198
 - vector space model, 103–104
- Word distributions
 - CPLSA, 424–425
 - LARA, 405
 - topics as, 335–340
- Word embedding in term clustering, 291–294
- Word-level ambiguity in NLP, 41
- Word relations, 251–252
- Word segmentation in NLP, 46
- Word sense disambiguation in NLP, 43
- Word valence scoring, 411
- Word vectors in text clustering, 278
- word2vec skip-gram, 293
- WordNet ontology, 294
- Zhai, ChengXiang, biography, 489
- Zipf's law
 - caching, 163
 - frequency analysis, 69–70

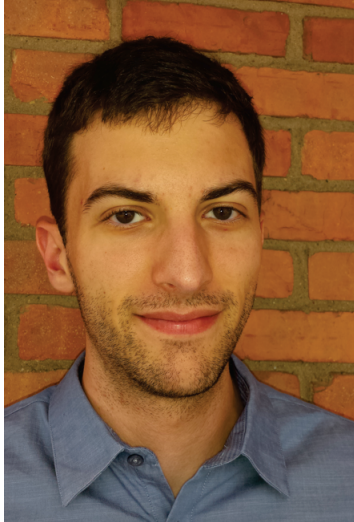
Authors' Biographies

ChengXiang Zhai



ChengXiang Zhai is a Professor of Computer Science and Willett Faculty Scholar at the University of Illinois at Urbana–Champaign, where he is also affiliated with the Graduate School of Library and Information Science, Institute for Genomic Biology, and Department of Statistics. He received a Ph.D. in Computer Science from Nanjing University in 1990, and a Ph.D. in Language and Information Technologies from Carnegie Mellon University in 2002. He worked at Clairvoyance Corp. as a Research Scientist and then Senior Research Scientist from 1997–2000. His research interests include information retrieval, text mining, natural language processing, machine learning, biomedical and health informatics, and intelligent education information systems. He has published over 200 research papers in major conferences and journals. He served as an Associate Editor for *Information Processing and Management*, as an Associate Editor of *ACM Transactions on Information Systems*, and on the editorial board of *Information Retrieval Journal*. He was a conference program co-chair of ACM CIKM 2004, NAACL HLT 2007, ACM SIGIR 2009, ECIR 2014, ICTIR 2015, and WWW 2015, and conference general co-chair for ACM CIKM 2016. He is an ACM Distinguished Scientist and a recipient of multiple awards, including the ACM SIGIR 2004 Best Paper Award, the ACM SIGIR 2014 Test of Time Paper Award, Alfred P. Sloan Research Fellowship, IBM Faculty Award, HP Innovation Research Program Award, Microsoft Beyond Search Research Award, and the Presidential Early Career Award for Scientists and Engineers (PECASE).

Sean Massung



Sean Massung is a Ph.D. candidate in computer science at the University of Illinois at Urbana-Champaign, where he also received both his B.S. and M.S. degrees. He is a co-founder of MeTA and uses it in all of his research. He has been instructor for CS 225: Data Structures and Programming Principles, CS 410: Text Information Systems, and CS 591txt: Text Mining Seminar. He is included in the 2014 List of Teachers Ranked as Excellent at the University of Illinois and has received an Outstanding Teaching Assistant Award and CS@Illinois Outstanding Research Project Award. He has given talks at Jump Labs Champaign and at UIUC for Data and Information Systems Seminar, Intro to Big Data, and Teaching Assistant Seminar. His research interests include text mining applications in information retrieval, natural language processing, and education.