
Sparse stochastic inference for latent Dirichlet allocation

Abstract

We present a hybrid algorithm for Bayesian topic models that combines the efficiency of sparse Gibbs sampling with the scalability of online stochastic inference. We used our algorithm to analyze a corpus of 1.2 million books (33 billion words) with thousands of topics. Our approach reduces the bias of variational inference and generalizes to many Bayesian hidden-variable models.

1. Electronic Submission

As in the past few years, ICML will rely exclusively on electronic formats for submission and review.

1.1. Templates for Papers

Electronic templates for producing papers for submission are available for L^AT_EX and Microsoft Word. Templates are accessible on the World Wide Web at: <http://icml.cc/2012/>

Send questions about these electronic templates to program@icml.cc.

The formatting instructions below will be enforced for initial submissions and camera-ready copies.

- The maximum paper length is 8 pages.
- Do not alter the style template; in particular, do not compress the paper format by reducing the vertical spaces.
- Do not include author information or acknowledgments in your initial submission.
- Place figure captions *under* the figure (and omit titles from inside the graphic file itself). Place table captions *over* the table.
- References must include page numbers whenever possible and be as complete as possible. Place multiple citations in chronological order.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

Please see below for details on each of these items.

1.2. Submitting Papers

Submission to ICML 2012 will be entirely electronic, via a web site (not email). The URL and information about the submission process are available on the conference web site at

<http://icml.cc/2012/>

Paper Deadline: The deadline for paper submission to ICML 2012 is Friday, February 24, 2012, at 23:59 Universal Time (3:59 Pacific Daylight Time). If your full submission does not reach us by this date, it will not be considered for publication. There is no separate abstract submission.

Anonymous Submission: To facilitate blind review, no identifying author information should appear on the title page or in the paper itself. Section 2.3 will explain the details of how to format this.

Simultaneous Submission: ICML will not accept any paper which, at the time of submission, is under review for another conference or has already been published. This policy also applies to papers that overlap substantially in technical content with conference papers under review or previously published. ICML submissions must not be submitted to other conferences during ICML's review period. Authors may submit to ICML substantially different versions of journal papers that are currently under review by the journal, but not yet accepted at the time of submission. Informal publications, such as technical reports or papers in workshop proceedings which do not appear in print, do not fall under these restrictions.

To ensure our ability to print submissions, authors must provide their manuscripts in **PDF** format. Furthermore, please make sure that files contain only Type-1 fonts (e.g., using the program **pdfonts** in linux or using File/DocumentProperties/Fonts in Acrobat). Other fonts (like Type-3) might come from graphics files imported into the document.

Authors using **Word** must convert their document to PDF. Most of the latest versions of Word have the facility to do this automatically. Submissions will not

be accepted in Word format or any format other than PDF. Really. We're not joking. Don't send Word.

$$M(s) < M(t) < |M| = m$$

$$y'' = c\{f[y', y(x)] + g(x)\}$$

$$\{f[g(x)]\}' = f'[g(x)]g'(x)$$

$$(f(x)g(x))'' = f''(x)g(x) + 2f'(x)g'(x) + f(x)g''(x)$$

$$\{f(g(x))\}' = f'(g(x))g'(x)$$

$$= f'(g(x))g'(x)$$

$$x^2 \qquad a_n \qquad x_i^n$$

$$x^{2n} \qquad x_{2y} \qquad A_{i,j,k}^{-n!2}$$

$$x^{y^2} \qquad x^{y_1} \qquad A_{j_{n,m}^{2n}}^2$$

$$\frac{1}{x+y} \qquad \frac{a^2-b^2}{a+b} = a-b$$

$$\sqrt[n]{\frac{x^n-y^n}{1+u^{2n}}} \qquad \frac{\frac{a}{x-y}+\frac{b}{x+y}}{1+\frac{a-b}{a+b}}$$

$$\sqrt{x^2+y^2+2xy}=x+y \qquad \sqrt[3]{-1+\sqrt{q^2+p^3}}$$

$$\sum_{i=1}^n \int_a^b$$

$$\int\limits_{x=0}^{x=1} 2\sum_{i=1}^n a_i \int_a^b f_i(x)g_i(x)\, \mathrm{d} x \\ a_0,a_1,\ldots,a_n \qquad a_0+a_1+\cdots+a_n$$

These are given by Cardan's formula as

$$y_1 = u + v$$

$$y_2 = -\frac{u+v}{2} + \frac{i}{2}\sqrt{3}(u-v)$$

$$y_3 = -\frac{u+v}{2} - \frac{i}{2}\sqrt{3}(u-v)$$

where

$$u = \sqrt[3]{-q + \sqrt{q^2 + p^3}}$$

$$v = \sqrt[3]{-q - \sqrt{q^2 + p^3}}$$

Each of the measurements $x_1 < x_2 < \cdots < x_r$, occurs p_1, p_2, \dots, p_r times. The mean value and standard deviation are then

$$x = \frac{1}{n} \sum_1^r p_i x_i$$

$$s = \sqrt{\frac{1}{n} \sum_1^r p_i (x_i - x)^2}$$

where $n = p_1 + p_2 + \cdots + p_r$.

Exercise 5.5: Although this equation looks very complicated, it should not present any great difficulties:

$$\int \frac{\sqrt{(ax+b)^3}}{x} \, dx =$$

$$\frac{2\sqrt{(ax+b)^3}}{3} + 2b\sqrt{ax+b} + b^2 \int \frac{dx}{x\sqrt{ax+b}}$$

The same applies to $\int_{-1}^8 (dx/\sqrt[3]{x}) = \frac{3}{2}(8^{\frac{2}{3}} + 1^{\frac{2}{3}}) = 15/2$.

Logistic Regression: Assume samples are $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$, $y^{(i)} \in \{1, 2, \dots, k\}$, θ represents a vector of all parameters, θ_j represents parameter corresponding to feature j . Each sample $x^{(i)}$'s probability is expressed as a sigmoid function. The probability that $x^{(i)}$ is assigned label 1 given θ is equal to

$$\begin{aligned} p(y^{(i)} = 1 | x^{(i)}; \theta) &= h_{\theta}(x) \\ &= \frac{1}{1 + \exp(-\theta^T x)} \\ &= \sigma(\theta^T x) \end{aligned} \quad (1)$$

So the objective function is

$$\begin{aligned} J(\theta) &= - \sum_i^m \{y^{(i)} \log(h_{\theta}(x^{(i)})) \\ &\quad + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\} \end{aligned} \quad (2)$$

The gradient is

$$\nabla_j J(\theta) = \sum_{i=1}^D x_j^{(i)} (h_{\theta}(x^{(i)}) - y^{(i)}) \quad (3)$$

Those who use \LaTeX to format their accepted papers need to pay close attention to the typefaces used. Specifically, when producing the PDF by first converting the dvi output of \LaTeX to Postscript the default behavior is to use non-scalable Type-3 PostScript bitmap fonts to represent the standard \LaTeX fonts. The resulting document is difficult to read in electronic form; the type appears fuzzy. To avoid this problem, dvips must be instructed to use an alternative font map. This can be achieved with something like the following commands:

dvips -Ppdf -tletter -G0 -o paper.ps paper.dvi ps2pdf paper.ps

Note that it is a zero following the “-G”. This tells dvips to use the config.pdf file (and this file refers to a better font mapping).

Another alternative is to use the **pdflatex** program instead of straight **L^AT_EX**. This program avoids the Type-3 font problem, however you must ensure that all of the fonts are embedded (use **pdffonts**). If they are not, you need to configure **pdflatex** to use a font map file that specifies that the fonts be embedded. Also you should ensure that images are not downsampled or otherwise compressed in a lossy way.

Note that the 2012 style files use the **hyperref** package to make clickable links in documents. If this causes problems for you, add **nohyperref** as one of the options to the **icml2012** usepackage statement.

1.3. Reacting to Reviews

We will continue the ICML tradition in which the authors are given the option of providing a short reaction to the initial reviews. These reactions $z = 2a + 3y$ will be taken into account in the discussion among the reviewers and area chairs.

$$2 \sum_{i=1}^n a_i \int_a^b f_i(x) g_i(x) dx \quad (4)$$

1.4. Submitting Final Camera-Ready Copy

The final versions of papers accepted for publication should follow the same format and naming convention as initial submissions, except of course that the normal author information (names and affiliations) should be given. See Section 2.3.2 for details of how to format this.

The footnote, “Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.” must be modified to “Appearing in *Proceedings of the 29th International Conference on Machine Learning*, Edinburgh, Scotland, UK, 2012. Copyright 2012 by the author(s)/owner(s).”

For those using the **L^AT_EX** style file, simply change `\usepackage{icml2012}` to

`\usepackage[accepted]{icml2012}`

Authors using **Word** must edit the footnote on the first page of the document themselves.

Camera-ready copies should have the title of the paper as running head on each page except the first one. The running title consists of a single line centered above a horizontal rule which is 1 point thick. The running head should be centered, bold and in 9 point type. The rule should be 10 points above the main text. For those using the **L^AT_EX** style file, the original title is automatically set as running head using the **fancyhdr** package which is included in the ICML 2012 style file

package. In case that the original title exceeds the size restrictions, a shorter form can be supplied by using

`\icmltitlerunning{...}`

just before `\begin{document}`. Authors using **Word** must edit the header of the document themselves.

2. Format of the Paper

All submissions must follow the same format to ensure the printer can reproduce them without problems and to let readers more easily find the information that they desire.

2.1. Length and Dimensions

Papers must not exceed eight (8) pages, including all figures, tables, references, and appendices. Any submission that exceeds this page limit or that diverges significantly from the format specified herein will be rejected without review.

The text of the paper should be formatted in two columns, with an overall width of 6.75 inches, height of 9.0 inches, and 0.25 inches between the columns. The left margin should be 0.75 inches and the top margin 1.0 inch (2.54 cm). The right and bottom margins will depend on whether you print on US letter or A4 paper, but all final versions must be produced for US letter size.

The paper body should be set in 10 point type with a vertical spacing of 11 points. Please use Times Roman typeface throughout the text.

2.2. Title

The paper title should be set in 14 point bold type and centered between two horizontal rules that are 1 point thick, with 1.0 inch between the top rule and the top edge of the page. Capitalize the first letter of content words and put the rest of the title in lower case.

2.3. Author Information for Submission

To facilitate blind review, author information must not appear. If you are using **L^AT_EX** and the **icml2012.sty** file, you may use `\icmlauthor{...}` to specify authors. The author information will simply not be printed until **accepted** is an argument to the style file. Submissions that include the author information will not be reviewed.

2.3.1. SELF-CITATIONS

If you are citing published papers for which you are an author, refer to yourself in the third person. In particular, do not use phrases that reveal your identity (e.g., “in previous work (Langley, 2000), we have shown ...”).

Do not anonymize citations in the reference section by removing or blacking out author names. The only exception are manuscripts that are not yet published (e.g. under submission). If you choose to refer to such unpublished manuscripts (Author, 2011), anonymized copies have to be submitted as Supplementary Material via CMT. However, keep in mind that an ICML paper should be self contained and should contain sufficient detail for the reviewers to evaluate the work. In particular, reviewers are not required to look at the Supplementary Material when writing their review.

2.3.2. CAMERA-READY AUTHOR INFORMATION

If a paper is accepted, a final camera-ready copy must be prepared. For camera-ready papers, author information should start 0.3 inches below the bottom rule surrounding the title. The authors’ names should appear in 10 point bold type, electronic mail addresses in 10 point small capitals, and physical addresses in ordinary 10 point type. Each author’s name should be flush left, whereas the email address should be flush right on the same line. The author’s physical address should appear flush left on the ensuing line, on a single line if possible. If successive authors have the same affiliation, then give their physical address only once.

A sample file (in PDF) with author names is included in the ICML2012 style file package.

2.4. Abstract

The paper abstract should begin in the left column, 0.4 inches below the final address. The heading ‘Abstract’ should be centered, bold, and in 11 point type. The abstract body should use 10 point type, with a vertical spacing of 11 points, and should be indented 0.25 inches more than normal on left-hand and right-hand margins. Insert 0.4 inches of blank space after the body. Keep your abstract brief and self-contained, limiting it to one paragraph and no more than six or seven sentences.

2.5. Partitioning the Text

You should organize your paper into sections and paragraphs to help readers place a structure on the material and understand its contributions.

2.5.1. SECTIONS AND SUBSECTIONS

Section headings should be numbered, flush left, and set in 11 pt bold type with the content words capitalized. Leave 0.25 inches of space before the heading and 0.15 inches after the heading.

Similarly, subsection headings should be numbered, flush left, and set in 10 pt bold type with the content words capitalized. Leave 0.2 inches of space before the heading and 0.13 inches afterward.

Finally, subsubsection headings should be numbered, flush left, and set in 10 pt small caps with the content words capitalized. Leave 0.18 inches of space before the heading and 0.1 inches after the heading.

Please use no more than three levels of headings.

2.5.2. PARAGRAPHS AND FOOTNOTES

Within each section or subsection, you should further partition the paper into paragraphs. Do not indent the first line of a given paragraph, but insert a blank line between succeeding ones.

You can use footnotes¹ to provide readers with additional information about a topic without interrupting the flow of the paper. Indicate footnotes with a number in the text where the point is most relevant. Place the footnote in 9 point type at the bottom of the column in which it appears. Precede the first footnote in a column with a horizontal rule of 0.8 inches.²

2.6. Figures

You may want to include figures in the paper to help readers visualize your approach and your results. Such artwork should be centered, legible, and separated from the text. Lines should be dark and at least 0.5 points thick for purposes of reproduction, and text should not appear on a gray background.

Label all distinct components of each figure. If the figure takes the form of a graph, then give a name for each axis and include a legend that briefly describes each curve. Do not include a title inside the figure; instead, the caption should serve this function.

Number figures sequentially, placing the figure number and caption *after* the graphics, with at least 0.1 inches of space before the caption and 0.1 inches after it, as in Figure 1. The figure caption should be set in 9 point

¹For the sake of readability, footnotes should be complete sentences.

²Multiple footnotes can appear in each column, in the same order as they appear in the text, but spread them across columns and pages if possible.

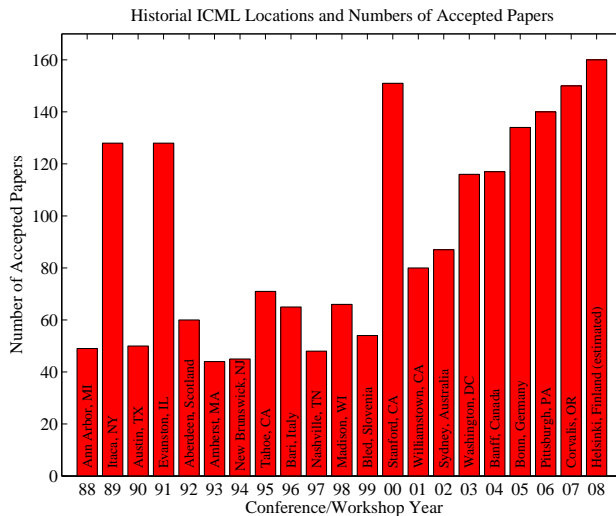


Figure 1. Historical locations and number of accepted papers for International Machine Learning Conferences (ICML 1993 – ICML 2008) and International Workshops on Machine Learning (ML 1988 – ML 1992). At the time this figure was produced, the number of accepted papers for ICML 2008 was unknown and instead estimated.

Algorithm 1 Bubble Sort

Input: data x_i , size m
repeat
 Initialize $noChange = true$.
 for $i = 1$ **to** $m - 1$ **do**
 if $x_i > x_{i+1}$ **then**
 Swap x_i and x_{i+1}
 $noChange = false$
 end if
 end for
until $noChange$ is $true$

Table 1. Classification accuracies for naive Bayes and flexible Bayes on various data sets.

DATA SET	NAIVE	FLEXIBLE	BETTER?
BREAST	95.9± 0.2	96.7± 0.2	✓
CLEVELAND	83.3± 0.6	80.0± 0.6	×
GLASS2	61.9± 1.4	83.8± 0.7	✓
CREDIT	74.8± 0.5	78.3± 0.6	
HORSE	73.3± 0.9	69.7± 1.0	×
META	67.1± 0.6	76.5± 0.5	✓
PIMA	75.1± 0.6	73.9± 0.5	
VEHICLE	44.9± 0.6	61.5± 0.4	✓

type and centered unless it runs two or more lines, in which case it should be flush left. You may float figures to the top or bottom of a column, and you may set wide figures across both columns (use the environment `figure*` in L^AT_EX), but always place two-column figures at the top or bottom of the page.

2.7. Algorithms

If you are using L^AT_EX, please use the “algorithm” and “algorithmic” environments to format pseudocode. These require the corresponding stylefiles, `algorithm.sty` and `algorithmic.sty`, which are supplied with this package. Algorithm 2 shows an example.

2.8. Tables

You may also want to include tables that summarize material. Like figures, these should be centered, legible, and numbered consecutively. However, place the title *above* the table with at least 0.1 inches of space before the title and the same after it, as in Table 1. The table title should be set in 9 point type and centered unless it runs two or more lines, in which case it should be flush left.

Tables contain textual material that can be typeset, as contrasted with figures, which contain graphical material that must be drawn. Specify the contents of each row and column in the table’s topmost row. Again, you may float tables to a column’s top or bottom, and set wide tables across both columns, but place two-

Algorithm 2 Bubble Sort listing all authors to a publication in an earlier reference

Step1 (data x_i , size m)

```

1: repeat
2:   Initialize  $noChange = true$ .
3:   for  $i = 1$  to  $m - 1$  do
4:     if  $x_i > x_{i+1}$  then
5:       Swap  $x_i$  and  $x_{i+1}$  Citations within the text
        should include the authors' last names and
6:        $noChange = false$ 
7:     end if
8:   end for
9: until  $noChange$  is true

```

Step2 (data x_i , size m)

```

1: repeat
2:   Initialize  $noChange = true$ .
3:   for  $i = 1$  to  $m - 1$  do
4:     if  $x_i > x_{i+1}$  then
5:       Swap  $x_i$  and  $x_{i+1}$ 
6:        $noChange = false$ 
7:     end if
8:   end for
9: until  $noChange$  is true

```

column tables at the top or bottom of the page.

2.9. Citations and References

Please use APA reference format regardless of your formatter or word processor. If you rely on the L^AT_EX bibliographic facility, use `natbib.sty` and `icml2012.bst` included in the style-file package to obtain this format.

Citations within the text should include the authors' last names and year. If the authors' names are included in the sentence, place only the year in parentheses, for example when referencing Arthur Samuel's pioneering work (1959). Otherwise place the entire reference in parentheses with the authors and year separated by a comma (Samuel, 1959). List multiple references separated by semicolons (Kearns, 1989; Samuel, 1959; Mitchell, 1980). Use the 'et al.' construct only for citations with three or more authors or after listing all authors to a publication in an earlier reference (Michalski et al., 1983).

Authors should cite their own work in the third person in the initial version of their paper submitted for blind review. Please refer to Section 2.3 for detailed instructions on how to cite your own papers.

Use an unnumbered first-level section heading for the references, and use a hanging indent style, with the first line of the reference flush against the left mar-

gin and subsequent lines indented by 10 points. The references at the end of this document give examples for journal articles (Samuel, 1959), conference publications (Langley, 2000), book chapters (Newell & Rosenbloom, 1981), books (Duda et al., 2000), edited volumes (Michalski et al., 1983), technical reports (Mitchell, 1980), and dissertations (Kearns, 1989).

Alphabetize references by the surnames of the first authors, with single author entries preceding multiple author entries. Order references for the same authors by year of publication, with the earliest first. Make sure that each reference includes all relevant information (e.g., page numbers).

2.10. Software and Data

We strongly encourage the publication of software and data with the camera-ready version of the paper whenever appropriate. This can be done by including a URL in the camera-ready copy. However, do not include URLs that reveal your institution or identity in your submission for review. Instead, provide an anonymous URL or upload the material as "Supplementary Material" into the CMT reviewing system. Note that reviewers are not required to look at this material when writing their review.

Acknowledgments

Do not include acknowledgements in the initial version of the paper submitted for blind review.

If a paper is accepted, the final camera-ready version can (and probably should) include acknowledgements. In this case, please place such acknowledgements in an unnumbered section at the end of the paper. Typically, this will include thanks to reviewers who gave useful comments, to colleagues who contributed to the ideas, and to funding agencies and corporate sponsors that provided financial support.

References

- Author, N. N. Suppressed for anonymity, 2011.
- Duda, R. O., Hart, P. E., and Stork, D. G. *Pattern Classification*. John Wiley and Sons, 2nd edition, 2000.
- Kearns, M. J. *Computational Complexity of Machine Learning*. PhD thesis, Department of Computer Science, Harvard University, 1989.
- Langley, P. Crafting papers on machine learning. In Langley, Pat (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML*

660	2000), pp. 1207–1216, Stanford, CA, 2000. Morgan	715
661	Kaufmann.	716
662		717
663	Michalski, R. S., Carbonell, J. G., and Mitchell, T. M.	718
664	(eds.). <i>Machine Learning: An Artificial Intelligence</i>	719
665	<i>Approach, Vol. I</i> . Tioga, Palo Alto, CA, 1983.	720
666		721
667	Mitchell, T. M. The need for biases in learning gener-	722
668	alizations. Technical report, Computer Science De-	723
669	partment, Rutgers University, New Brunswick, MA,	724
670	1980.	725
671		726
672	Newell, A. and Rosenbloom, P. S. Mechanisms of skill	727
673	acquisition and the law of practice. In Anderson,	728
674	J. R. (ed.), <i>Cognitive Skills and Their Acquisition</i> ,	729
675	chapter 1, pp. 1–51. Lawrence Erlbaum Associates,	730
676	Inc., Hillsdale, NJ, 1981.	731
677		732
678	Samuel, A. L. Some studies in machine learning using	733
679	the game of checkers. <i>IBM Journal of Research and</i>	734
680	<i>Development</i> , 3(3):211–229, 1959.	735
681		736
682		737
683		738
684		739
685		740
686		741
687		742
688		743
689		744
690		745
691		746
692		747
693		748
694		749
695		750
696		751
697		752
698		753
699		754
700		755
701		756
702		757
703		758
704		759
705		760
706		761
707		762
708		763
709		764
710		765
711		766
712		767
713		768
714		769