

Self-Tuning for Sentiment Analysis in Scientific Citation

Huy Tu

North Carolina State University
Raleigh, North Carolina
hqtu@ncsu.edu

ABSTRACT

Sentiment analysis of citations within scientific papers have been studied to more appropriately indicating the quality of published papers in comparison to the outdated quantitative evaluation approach basing on the citations' frequency. Approximating the popularity, the context, and the impact of the published research are commonly referred to as bibliometric measurements. In this study, we shall review state of the art sentiment analysis algorithms that helps incorporated that qualitative aspect to the bibliometrics. Moreover, we will focus on the improvement of result analysis of those state of the art algorithms by proposing empirical software science approaches such as ensemble methods and parameter tuning to.

Keywords - Citation Sentiment Analysis, Classification; Differential Evolution, Parameters Tuning; Clustering, Cluster Ensemble.

1 INTRODUCTION

Bibliographic citations have kept the research community alive. Iorio et al, Nuzzolese et al, and Peroni et al have expressed the importance of bibliographic citations as they are tools for:

- *disseminating research* - with references as directional edge to point to motivational background, related work, data sources, and methods coming from various publishing platforms.
- *exploring research* - the network of citations (papers or authors or journals are presented by the nodes and the edges representing the referencing act) can be utilized as readable data for searching, visualizing, filtering, and aggregating for interesting problems such as knowledge transfers across domains, scientific collaboration, and the scholarly disciplines map.
- *evaluating research* - characterizing and quantifying the importance and impact of the work/research through the nature of the purposes those citations and not only the mere existence of those papers.

In this paper, we will focus only at the third important aspect of bibliographic citations, evaluating research. Evaluating research's impact should be optimized for fairness and objectiveness through both quantitative and qualitative evaluations. Only considering one without the other aspect of the evaluation would be misinformed, narrow-minded, and problematic.

- For quantitative aspect without qualitative aspect, the nature of citation itself would not be considered and all the citations would be taken in as they carry equivalent weights when they should not. For example, the referencing act of crediting should be counted differently for the referencing act of criticizing.
- For qualitative aspect without quantitative aspect, the interest of other researcher in it and the popularity of the work

[] would be missing. Moreover, there would be no statistics for significance and effect insights from the qualitative evaluation.

The old-fashioned quantitative citation bibliometrics include Hirsh-index [], the g-index [], and PageRank []. The modern qualitative NLP-based bibliometrics include supervised learning of polarity and purpose of citation classification.

Citation purpose classification is the studying of the motive behind citing that work or research. In 1977, Spiefl-Rosing et al. [] formerly suggested 13 categories for citation purpose. In 2006, Teufel et al. [] adopted 12 categories from Spiegel-Rosing's work that can be appropriately grouped to four types: weakness, contrast (4 categories), positive (6 categories), and neutral. In 2013, from those previous work, Abu-Jbara et al. [] condensed the previous's quantity of categories down to only six categories of: criticizing, comparison, use, substantiating, basis, and neutral(other).

Citation polarity classification through sentiment analysis aims to determine opinions, emotions, and attitudes of the specific granularity of the text region. It can be conducted at three levels of document-level, sentence-level, and context-level []. We will focus more on citation polarity classification on context level for this work. The document-level and sentence-level of sentiment analysis have been incomplete and defective due to how citation region can be varied, how multiple citations can be included in one sentence while having different sentiment tags, and how the sentiment can be hidden in the context level. For example,

2 Related Work

The work of Och et al (2004) is perhaps the best-known study of new features and their impact on translation quality. However, it had a few shortcomings. First, it used the features for reranking n -best lists of translations, rather than for decoding or forest reranking (Huang, 2008). Second, it attempted to incorporate syntax by applying off-the-shelf part-of-speech taggers and parsers to MT output, a task these tools were never designed for. By contrast, we incorporate features directly into hierarchical and syntax-based decoders.

A third difficulty with Och et al.'s study was that it used MERT, which is not an ideal vehicle for feature exploration because it is observed not to perform well with large feature sets. Others have in-

Figure 1: A sample black and white graphic that has been resized with the includegraphics command.

In the figure 1 above, the citation of Och et al from "Smorgasbord of Features for Statistical Machine Translation" paper firstly

included 'best known study', making this citation positive. However, if we read the following sentences after this citation, we find that this citation is being referred in the next 5 sentences anaphorically like the word 'it' and other such phrases which criticizing the original referenced work negatively which makes this citation sentiment difficult to judge.

Specifically, there are still many difficulties and open research opportunities within the existing models for citation sentiment analysis:

- Sentiment is often subtle or even hidden in citation within research community. Negative polarity is often embedded contrastively.
- Sentiment of the citation within sentences are often neutral.
- Diverse variation of sentiment lexicon or sentiment carrying scientific or technical terms/phrases. Some are longer in length (such as "state of the art") which suggests that considering higher n-grams would be useful.
- From a single clause level to multiple paragraphs level, the region of influence of citations differ.

2 THE BODY OF THE PAPER

Typically, the body of a paper is organized into a hierarchical structure, with numbered or unnumbered headings for sections, subsections, sub-subsections, and even smaller sections. The command `\section` that precedes this paragraph is part of such a hierarchy.¹ \LaTeX handles the numbering and placement of these headings for you, when you use the appropriate heading commands around the titles of the headings. If you want a sub-subsection or smaller part to be unnumbered in your output, simply append an asterisk to the command name. Examples of both numbered and unnumbered headings will appear throughout the balance of this sample document.

Because the entire article is contained in the **document** environment, you can indicate the start of a new paragraph with a blank line in your input file; that is why this sentence forms a separate paragraph.

2.1 Type Changes and *Special* Characters

We have already seen several typeface changes in this sample. You can indicate italicized words or phrases in your text with the command `\textit`; boldening with the command `\textbf` and typewriter-style (for instance, for computer code) with `\texttt`. But remember, you do not have to indicate typestyle changes when such changes are part of the *structural* elements of your article; for instance, the heading of this subsection will be in a sans serif² typeface, but that is handled by the document class file. Take care with the use of³ the curly braces in typeface changes; they mark the beginning and end of the text that is to be in the different typeface.

You can use whatever symbols, accented characters, or non-English characters you need anywhere in your document; you can find a complete list of what is available in the *\LaTeX User's Guide* [26].

¹This is a footnote.

²Another footnote here. Let's make this a rather long one to see how it looks.

³Another footnote.

2.2 Math Equations

You may want to display math equations in three distinct styles: inline, numbered or non-numbered display. Each of the three are discussed in the next sections.

2.2.1 Inline (In-text) Equations. A formula that appears in the running text is called an inline or in-text formula. It is produced by the **math** environment, which can be invoked with the usual `\begin . . . \end` construction or with the short form `$. . . $`. You can use any of the symbols and structures, from α to ω , available in \LaTeX [26]; this section will simply show a few examples of in-text equations in context. Notice how this equation: $\lim_{n \rightarrow \infty} x = 0$, set here in in-line math style, looks slightly different when set in display style. (See next section).

2.2.2 Display Equations. A numbered display equation—one set off by vertical space from the text and centered horizontally—is produced by the **equation** environment. An unnumbered display equation is produced by the **displaymath** environment.

Again, in either environment, you can use any of the symbols and structures available in \LaTeX ; this section will just give a couple of examples of display equations in context. First, consider the equation, shown as an inline equation above:

$$\lim_{n \rightarrow \infty} x = 0 \quad (1)$$

Notice how it is formatted somewhat differently in the **displaymath** environment. Now, we'll enter an unnumbered equation:

$$\sum_{i=0}^{\infty} x + 1$$

and follow it with another numbered equation:

$$\sum_{i=0}^{\infty} x_i = \int_0^{\pi+2} f \quad (2)$$

just to demonstrate \LaTeX 's able handling of numbering.

2.3 Citations

Citations to articles [6–8, 19], conference proceedings [8] or maybe books [26, 34] listed in the Bibliography section of your article will occur throughout the text of your article. You should use BibTeX to automatically produce this bibliography; you simply need to insert one of several citation commands with a key of the item cited in the proper location in the `.tex` file [26]. The key is a short reference you invent to uniquely identify each work; in this sample document, the key is the first author's surname and a word from the title. This identifying key is included with each item in the `.bib` file for your article.

The details of the construction of the `.bib` file are beyond the scope of this sample document, but more information can be found in the *Author's Guide*, and exhaustive details in the *\LaTeX User's Guide* by Lammport [26].

This article shows only the plainest form of the citation command, using `\cite`.

Some examples. A paginated journal article [2], an enumerated journal article [11], a reference to an entire issue [10], a monograph (whole book) [25], a monograph/whole book in a series (see 2a in spec. document) [18], a divisible-book such as an anthology or

Table 1: Citation Purposes

Category	Description
Critique	The citing sentence mentions the weakness/strengths of the cited approach, negatively/positively criticizes the cited approach, negatively/positively evaluates the cited source.
Comparison	The citing sentence compares or contrasts the work in the cited paper to the author’s work.
Use	The citing sentence uses the method, idea or tool of the cited paper.
Substantiating	The results, claims of the citing work substantiate, verify the cited paper and support each other.
Basis	The author uses the cited work as starting point or motivation and extends on the cited work
Neutral (Other)	The sentence is a neutral description of the cited work or if it doesn’t come under any of the above categories.

Table 2: Frequency of Special Characters

Non-English or Math	Frequency	Comments
Ø	1 in 1,000	For Swedish names
π	1 in 5	Common in math
\$	4 in 5	Used in business
Ψ ₁ ²	1 in 40,000	Unexplained usage

compilation [13] followed by the same example, however we only output the series if the volume number is given [14] (so Editor00a’s series should NOT be present since it has no vol. no.), a chapter in a divisible book [37], a chapter in a divisible book in a series [12], a multi-volume work as book [24], an article in a proceedings (of a conference, symposium, workshop for example) (paginated proceedings article) [4], a proceedings article with all possible elements [36], an example of an enumerated proceedings article [16], an informally published work [17], a doctoral dissertation [9], a master’s thesis: [5], an online document / world wide web resource [1, 30, 38], a video game (Case 1) [29] and (Case 2) [28] and [27] and (Case 3) a patent [35], work accepted for publication [31], ‘YYYYb’-test for prolific author [32] and [33]. Other cites might contain ‘duplicate’ DOI and URLs (some SIAM articles) [23]. Boris / Barbara Beeton: multi-volume works as books [21] and [20].

A couple of citations with DOIs: [22, 23].
Online citations: [38–40].

2.4 Tables

Because tables cannot be split across pages, the best placement for them is typically the top of the page nearest their initial cite. To ensure this proper “floating” placement of tables, use the environment **table** to enclose the table’s contents and the table caption. The contents of the table itself must go in the **tabular** environment, to be aligned properly in rows and columns, with the desired horizontal and vertical rules. Again, detailed instructions on **tabular** material are found in the *L^AT_EX User’s Guide*.

Immediately following this sentence is the point at which Table 2 is included in the input file; compare the placement of the table here with the table in the printed output of this document.

To set a wider table, which takes up the whole width of the page’s live area, use the environment **table*** to enclose the table’s contents and the table caption. As with a single-column table, this wide table will “float” to a location deemed more desirable. Immediately following this sentence is the point at which Table 3 is included in



Figure 2: A sample black and white graphic.



Figure 3: A sample black and white graphic that has been resized with the `includegraphics` command.

the input file; again, it is instructive to compare the placement of the table here with the table in the printed output of this document. It is strongly recommended to use the package `booktabs` [15] and follow its main principles of typography with respect to tables:

- (1) Never, ever use vertical rules.
- (2) Never use double rules.

It is also a good idea not to overuse horizontal rules.

2.5 Figures

Like tables, figures cannot be split across pages; the best placement for them is typically the top or the bottom of the page nearest their initial cite. To ensure this proper “floating” placement of figures, use the environment **figure** to enclose the figure and its caption.

This sample document contains examples of .eps files to be displayable with L^AT_EX. If you work with pdfL^AT_EX, use files in the .pdf format. Note that most modern T_EX systems will convert .eps to .pdf for you on the fly. More details on each of these are found in the *Author’s Guide*.

As was the case with tables, you may want a figure that spans two columns. To do this, and still to ensure proper “floating” placement of tables, use the environment **figure*** to enclose the figure and its caption. And don’t forget to end the environment with **figure***, not **figure**!

Table 3: Some Typical Commands

Command	A Number	Comments
\author	100	Author
\table	300	For tables
\table*	400	For wider tables

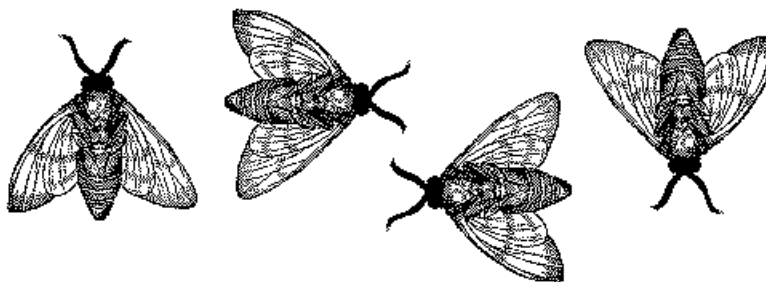


Figure 4: A sample black and white graphic that needs to span two columns of text.

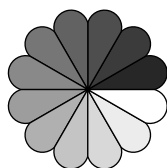


Figure 5: A sample black and white graphic that has been resized with the includegraphics command.

2.6 Theorem-like Constructs

Other common constructs that may occur in your article are the forms for logical constructs like theorems, axioms, corollaries and proofs. ACM uses two types of these constructs: theorem-like and definition-like.

Here is a theorem:

THEOREM 2.1. *Let f be continuous on $[a, b]$. If G is an antiderivative for f on $[a, b]$, then*

$$\int_a^b f(t) dt = G(b) - G(a).$$

Here is a definition:

Definition 2.2. If z is irrational, then by e^z we mean the unique number that has logarithm z :

$$\log e^z = z.$$

The pre-defined theorem-like constructs are **theorem**, **conjecture**, **proposition**, **lemma** and **corollary**. The pre-defined definition-like constructs are **example** and **definition**. You can add your own constructs using the *amsthm* interface [3]. The styles used in the `\theoremstyle` command are **acmplain** and **acmdefinition**.

Another construct is **proof**, for example,

PROOF. Suppose on the contrary there exists a real number L such that

$$\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = L.$$

Then

$$l = \lim_{x \rightarrow c} f(x) = \lim_{x \rightarrow c} \left[g(x) \cdot \frac{f(x)}{g(x)} \right] = \lim_{x \rightarrow c} g(x) \cdot \lim_{x \rightarrow c} \frac{f(x)}{g(x)} = 0 \cdot L = 0,$$

which contradicts our assumption that $l \neq 0$. \square

3 CONCLUSIONS

This paragraph will end the body of this sample document. Remember that you might still have Acknowledgments or Appendices; brief samples of these follow. There is still the Bibliography to deal with; and we will make a disclaimer about that here: with the exception of the reference to the \LaTeX book, the citations in this paper are to articles which have nothing to do with the present subject and are used as examples only.

A HEADINGS IN APPENDICES

The rules about hierarchical headings discussed above for the body of the article are different in the appendices. In the **appendix** environment, the command **section** is used to indicate the start of each Appendix, with alphabetic order designation (i.e., the first is A, the second B, etc.) and a title (if you include one). So, if you need hierarchical structure *within* an Appendix, start with **subsection** as the highest level. Here is an outline of the body of this document in Appendix-appropriate form:

A.1 Introduction

A.2 The Body of the Paper

A.2.1 Type Changes and Special Characters.

A.2.2 Math Equations.

Inline (In-text) Equations.

Display Equations.

A.2.3 Citations.

A.2.4 Tables.

A.2.5 Figures.

A.2.6 Theorem-like Constructs.

A Caveat for the \TeX Expert.

A.3 Conclusions

A.4 References

Generated by bibtex from your .bib file. Run latex, then bibtex, then latex twice (to resolve references) to create the .bbl file. Insert that .bbl file into the .tex source file and comment out the command \thebibliography.

B MORE HELP FOR THE HARDY

Of course, reading the source code is always useful. The file acmart.pdf contains both the user guide and the commented code.

ACKNOWLEDGMENTS

The authors would like to thank Dr. Yuhua Li for providing the MATLAB code of the BEPS method.

The authors would also like to thank the anonymous referees for their valuable comments and helpful suggestions. The work is supported by the National Natural Science Foundation of China under Grant No.: 61273304 and Young Scientists' Support Program (<http://www.nnsf.cn/youngscientists>).

REFERENCES

- [1] Rafal Ablamowicz and Bertfried Fauser. 2007. CLIFFORD: a Maple 11 Package for Clifford Algebra Computations, version 11. (2007). Retrieved February 28, 2008 from <http://math.tntech.edu/rafal/cliff11/index.html>
- [2] Patricia S. Abril and Robert Plant. 2007. The patent holder's dilemma: Buy, sell, or troll? *Commun. ACM* 50, 1 (Jan. 2007), 36–44. <https://doi.org/10.1145/1188913.1188915>
- [3] American Mathematical Society 2015. *Using the amsthm Package*. American Mathematical Society. <http://www.ctan.org/pkg/amsthm>.
- [4] Sten Andler. 1979. Predicate Path expressions. In *Proceedings of the 6th. ACM SIGACT-SIGPLAN symposium on Principles of Programming Languages (POPL '79)*. ACM Press, New York, NY, 226–236. <https://doi.org/10.1145/567752.567774>
- [5] David A. Anisi. 2003. *Optimal Motion Control of a Ground Vehicle*. Master's thesis. Royal Institute of Technology (KTH), Stockholm, Sweden.
- [6] Mic Bowman, Saumya K. Debray, and Larry L. Peterson. 1993. Reasoning About Naming Systems. *ACM Trans. Program. Lang. Syst.* 15, 5 (November 1993), 795–825. <https://doi.org/10.1145/161468.161471>
- [7] Johannes Braams. 1991. Babel, a Multilingual Style-Option System for Use with LaTeX's Standard Document Styles. *TUGboat* 12, 2 (June 1991), 291–301.
- [8] Malcolm Clark. 1991. Post Congress Tristesse. In *TeX90 Conference Proceedings*. TeX Users Group, 84–89.
- [9] Kenneth L. Clarkson. 1985. *Algorithms for Closest-Point Problems (Computational Geometry)*. Ph.D. Dissertation. Stanford University, Palo Alto, CA. UMI Order Number: AAT 8506171.
- [10] Jacques Cohen (Ed.). 1996. Special issue: Digital Libraries. *Commun. ACM* 39, 11 (Nov. 1996).
- [11] Sarah Cohen, Werner Nutt, and Yehoshua Sagie. 2007. Deciding equivalences among conjunctive aggregate queries. *J. ACM* 54, 2, Article 5 (April 2007), 50 pages. <https://doi.org/10.1145/1219092.1219093>
- [12] Bruce P. Douglass, David Harel, and Mark B. Trakhtenbrot. 1998. Statecarts in use: structured analysis and object-orientation. In *Lectures on Embedded Systems*, Grzegorz Rozenberg and Frits W. Vaandrager (Eds.). Lecture Notes in Computer Science, Vol. 1494. Springer-Verlag, London, 368–394. https://doi.org/10.1007/3-540-65193-4_29
- [13] Ian Editor (Ed.). 2007. *The title of book one* (1st. ed.). The name of the series one, Vol. 9. University of Chicago Press, Chicago. <https://doi.org/10.1007/3-540-09237-4>
- [14] Ian Editor (Ed.). 2008. *The title of book two* (2nd. ed.). University of Chicago Press, Chicago, Chapter 100. <https://doi.org/10.1007/3-540-09237-4>
- [15] Simon Fear. 2005. *Publication quality tables in \TeX* . <http://www.ctan.org/pkg/booktabs>.
- [16] Matthew Van Gundy, Davide Balzarotti, and Giovanni Vigna. 2007. Catch me, if you can: Evading network signatures with web-based polymorphic worms. In *Proceedings of the first USENIX workshop on Offensive Technologies (WOOT '07)*. USENIX Association, Berkley, CA, Article 7, 9 pages.
- [17] David Harel. 1978. *LOGICS of Programs: AXIOMATICS and DESCRIPTIVE POWER*. MIT Research Lab Technical Report TR-200. Massachusetts Institute of Technology, Cambridge, MA.
- [18] David Harel. 1979. *First-Order Dynamic Logic*. Lecture Notes in Computer Science, Vol. 68. Springer-Verlag, New York, NY. <https://doi.org/10.1007/3-540-09237-4>
- [19] Maurice Herlihy. 1993. A Methodology for Implementing Highly Concurrent Data Objects. *ACM Trans. Program. Lang. Syst.* 15, 5 (November 1993), 745–770. <https://doi.org/10.1145/161468.161469>
- [20] Lars Hörmander. 1985. *The analysis of linear partial differential operators. III. Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*, Vol. 275. Springer-Verlag, Berlin, Germany. viii+525 pages. Pseudodifferential operators.
- [21] Lars Hörmander. 1985. *The analysis of linear partial differential operators. IV. Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*, Vol. 275. Springer-Verlag, Berlin, Germany. vii+352 pages. Fourier integral operators.
- [22] IEEE 2004. IEEE TCSC Executive Committee. In *Proceedings of the IEEE International Conference on Web Services (ICWS '04)*. IEEE Computer Society, Washington, DC, USA, 21–22. <https://doi.org/10.1109/ICWS.2004.64>
- [23] Markus Kirschmer and John Voight. 2010. Algorithmic Enumeration of Ideal Classes for Quaternion Orders. *SIAM J. Comput.* 39, 5 (Jan. 2010), 1714–1747. <https://doi.org/10.1137/080734467>
- [24] Donald E. Knuth. 1997. *The Art of Computer Programming, Vol. 1: Fundamental Algorithms (3rd. ed.)*. Addison Wesley Longman Publishing Co., Inc.
- [25] David Kosior. 2001. *Understanding Policy-Based Networking* (2nd. ed.). Wiley, New York, NY.
- [26] Leslie Lamport. 1986. *\TeX : A Document Preparation System*. Addison-Wesley, Reading, MA.
- [27] Newton Lee. 2005. Interview with Bill Kinder: January 13, 2005. Video. *Comput. Entertain.* 3, 1, Article 4 (Jan.-March 2005). <https://doi.org/10.1145/1057270.1057278>
- [28] Dave Novak. 2003. Solder man. Video. In *ACM SIGGRAPH 2003 Video Review on Animation theater Program: Part I - Vol. 145 (July 27–27, 2003)*. ACM Press, New York, NY, 4. <https://doi.org/99.9999/woot07-S422>
- [29] Barack Obama. 2008. A more perfect union. Video. (5 March 2008). Retrieved March 21, 2008 from <http://video.google.com/videoplay?docid=6528042696351994555>
- [30] Poker-Edge.Com. 2006. Stats and Analysis. (March 2006). Retrieved June 7, 2006 from <http://www.poker-edge.com/stats.php>
- [31] Bernard Rous. 2008. The Enabling of Digital Libraries. *Digital Libraries* 12, 3, Article 5 (July 2008). To appear.
- [32] Mehdi Saeedi, Morteza Saheb Zamani, and Mehdi Sedighi. 2010. A library-based synthesis methodology for reversible logic. *Microelectron. J.* 41, 4 (April 2010), 185–194.
- [33] Mehdi Saeedi, Morteza Saheb Zamani, Mehdi Sedighi, and Zahra Sasanian. 2010. Synthesis of Reversible Circuit Using Cycle-Based Approach. *J. Emerg. Technol. Comput. Syst.* 6, 4 (Dec. 2010).
- [34] S.L. Salas and Einar Hille. 1978. *Calculus: One and Several Variable*. John Wiley and Sons, New York.
- [35] Joseph Scientist. 2009. The fountain of youth. (Aug. 2009). Patent No. 12345, Filed July 1st., 2008, Issued Aug. 9th., 2009.
- [36] Stan W. Smith. 2010. An experiment in bibliographic mark-up: Parsing metadata for XML export. In *Proceedings of the 3rd. annual workshop on Librarians and Computers (LAC '10)*, Reginald N. Smythe and Alexander Noble (Eds.), Vol. 3. Paparazzi Press, Milan Italy, 422–431. <https://doi.org/99.9999/woot07-S422>
- [37] Asad Z. Spector. 1990. Achieving application requirements. In *Distributed Systems* (2nd. ed.), Sape Mullender (Ed.). ACM Press, New York, NY, 19–33. <https://doi.org/10.1145/90417.90738>
- [38] Harry Thornburg. 2001. Introduction to Bayesian Statistics. (March 2001). Retrieved March 2, 2005 from <http://ccrma.stanford.edu/~jos/bayes/bayes.html>
- [39] TUG 2017. Institutional members of the TeX Users Group. (2017). Retrieved May 27, 2017 from <http://www.tug.org/institution.html>
- [40] Boris Veytsman. [n. d.]. acmart—Class for typesetting publications of ACM. ([n. d.]). Retrieved May 27, 2017 from <http://www.ctan.org/pkg/acmart>