

An Ontology of Privacy Law Derived from Published Articles Using Latent Dirichlet Allocation Probabilistic Topic Modeling

1 **Robert Sprague**

2 Department of Management & Marketing
3 University of Wyoming
4 Laramie, WY 82071
5 spraguer@uwyo.edu
6

Nicole Barberis

 IBM
 Laramie, WY
 nbarberis@yahoo.com

7 **Abstract**

8 This paper presents preliminary results from applying latent Dirichlet
9 allocation probabilistic topic modeling algorithms to a document collection
10 comprised of published law review articles from 1891 through 1970, and
11 sample publications from 1980, 1990, 2000, and 2010—all citing to
12 Warren’s and Brandeis’s *The Right to Privacy* and substantively discussing
13 privacy law issues. Our initial interpretation of the results reveals a few
14 important trends: early discussion of issues associated with the tort of
15 appropriation of name or likeness; fourth amendment government searches
16 and surveillance; and emerging trends in information privacy associated
17 with data collection, as well as the rise of the Internet. Once our initial data
18 collection has been completed, by filling out the years 1971 through 2012,
19 we hope to not only confirm these initial observations but also reveal
20 additional privacy law trends, developing an ontology of privacy law based
21 upon the topics revealed in published law review and journal articles.

22 23 **1 Introduction**

24 Privacy, being an evolutionary product of social development[3], has been a human need and
25 desire for millennia. Privacy law scholarship, in contrast, is a relatively recent phenomenon.
26 Within this recent profusion of scholarship lies a conundrum: there is no clear definition of
27 privacy[5]; there is not even consensus of what would constitute an adequate description.
28 Fundamental concepts associated with privacy have been identified and analyzed—for
29 example, seclusion, intimacy, surveillance, anonymity, and control of information. But, as
30 Helen Nissenbaum has noted, most calls for privacy arise from context, as well as advancing
31 technologies[4], meaning the legal system often has difficulty identifying and protecting
32 rights to privacy. Without a coherent construction of privacy principles shared by the
33 community of scholars, the legal discipline will never explicitly articulate those
34 principles[5].

35 This paper reports preliminary results from a research project aimed at identifying
36 fundamental privacy law principles derived from the writings of legal scholars and
37 commentators using probabilistic topic modeling. A latent Dirichlet allocation (LDA)
38 process, which identifies sets of terms that more tightly co-occur, is incorporated into the
39 topic modeling analysis to identify words most closely associated with each identified topic.
40 The LDA therefore provides insight into the context in which each identified topic occurs.

41 Most published law review articles that cite to Samuel Warren’s and Louis Brandeis’s

42 seminal article, *The Right to Privacy* [6] (some 3000 articles), are being converted to plain
43 text. *The Right to Privacy* was selected as the focal point of the document collection because
44 it is the original published scholarly call for a formal legal right to privacy in the United
45 States; hence, the vast majority of privacy law publications cites to it. Probabilistic topic
46 modeling using latent Dirichlet allocation is being applied to the document collection in time
47 slices to reveal the evolution of fundamental privacy law concepts expressed in the legal
48 literature published from 1890 through 2012. The ultimate goal of this project is to identify
49 the fundamental conceptual structure of privacy law in the United States as reflected by over
50 a century of published law review and journal articles.

52 **2 Methodology**

53 The initial document collection will be comprised of relevant published law review and
54 journal articles that cite to Warren’s and Brandeis’s *The Right to Privacy*, as identified
55 within the Westlaw and HeinOnline collections. All documents selected for the collection
56 are converted to plain text. In addition, all titles, author names, section headings,
57 footnotes/endnotes, and supplemental materials are removed in an effort to create a
58 collection limited to addressing substantive privacy law issues. At present, only
59 approximately 20% of the anticipated initial collection, representing privacy law articles
60 published up to and through 1970, has been converted. The current document collection has
61 been divided into the following time-slice corpora: 1891-1940 (which includes *The Right to*
62 *Privacy*), 1941-1950, 1951-1960, and 1961-1970. A cumulative corpus has also been created
63 for the time period 1891-1970. For this paper, additional partial corpora were created from
64 articles published in 1980, 1990, 2000, and 2010. Work is continuing to build complete
65 corpora for: 1971-1980, 1981-1990, 1991-2000, 2001-2010, and 2010-2012.

66 The LDA topic model algorithms are then applied to the corpora using MALLET[2]. The
67 critical MALLET output files used in this project include the following files: keys, weights,
68 words count, and composition.

70 **2.1 Results**

71 Identifying significant privacy law topics from the individual corpora can be approached
72 from different views of the MALLET results. The MALLET output files were parsed and
73 analyzed using R, which also generated visualizations of the data.

74 **2.1.1 Ubiquity Measure**

75 Our “Ubiquity Measure” is an attempt to visualize the frequency of occurrence of the term
76 “privacy” throughout the individual corpora. Within each time-slice corpus, every
77 occurrence of the term “privacy” was assigned an “AB” weight (where A = the normalized
78 weight of the topic in which the term occurs and B = the normalized weight of the term
79 within its topic). The AB weight is therefore the adjusted weight of the term depending on
80 the weight of the topic within which it occurs. Our “Ubiquity Measure” is the sum of all the
81 AB values within a particular time-slice corpus. In effect, the “Ubiquity Measure” represents
82 the degree of occurrence of the term privacy in each time-slice relative to all the other time
83 slices.

84 **2.1.2 Privacy Constellations**

85 Our “Privacy Constellations” reflect the normalized weight of the term privacy within each
86 topic in which it occurs within each time-slice corpus. This not only reflects the other terms
87 co-occurring with the term privacy within a topic, but also each term’s relative weight within
88 the topic. This allows one to see the context in which the term privacy was used by
89 published authors within each time-slice corpus.

90 **2.1.3 Treemaps**

91 We created Treemaps for each time-slice corpus reflecting the weight of the topics in the
92 corpus and labeled each area with the first term in the topic’s cluster because the
93 “discovered” topic is mostly about that first word. This is a good way to see which ideas
94 were important in an era, reflected by their relative weight.

95 **2.1.4 Topics and Terms**

96 The topics and associated terms identified for each time slice are reflected in a 20x20 matrix
97 built from MALLETT's "keys" output file. However, these matrices can be somewhat
98 cumbersome to read and interpret. We believe our Ubiquity Measure, Privacy Constellations,
99 and Treemaps offer more helpful visualizations of our data for interpreting the results.

100

101 **2.2 Initial Interpretations of the Preliminary Results**

102 Due to the page limitation for this proposal—targeting the Application theme of the NIPS
103 Topic Models workshop—visualizations of our initial results are not included. If this
104 proposal is accepted into the workshop, the visualizations can comprise a significant portion
105 of this project's presentation.

106 We can make a few generalizations from the preliminary results. Looking at the Treemaps, it
107 is not surprising to find the term "privacy" as one of the (normalized) heaviest-weighted
108 "top" terms. The Treemaps also provide a few additional insights: for example, the term
109 property was a heavily weighted term through 1950, but then subsequently drops out of the
110 "top" terms. Meanwhile, the terms government, fourth, and amendment make sizable
111 appearances in the 1961-1970 corpus, implying much greater attention to fourth amendment
112 privacy rights related to government searches.¹

113 In the 1891-1940 corpus, the term "privacy" is closely associated with the terms "public,"
114 "publication," "picture," "person," "interest," "invasion," "news," "life," "advertising," and
115 "photograph." This clustering reflects a focus in the early development of privacy law on the
116 ability of individuals to control their images and likenesses—fundamentally, the tort of
117 appropriation of name or likeness. As the document collection is completed beyond 1970 it
118 will be interesting to observe how prevalent this topic will remain in comparison to other
119 "privacy" topics.

120 The term "information" appears quite frequently within the corpus, although its most
121 frequent appearance is in a very low-weighted topic and most frequently co-occurs with the
122 terms "data," "credit," "computer," "personal," "privacy," "bureau," "access," "system," and
123 "files" (we can infer "credit" and "bureau" refer to credit bureaus because those two terms
124 most tightly co-occur also in the same topic). These terms can be associated with the
125 growing computerized collection of personal information that began in the 1960s. The fact
126 that many of these terms, particularly "information," appear frequently in the corpus, but
127 most frequently in a very low-weighted topic, indicates that this privacy issue was a late
128 blooming topic, at least for this corpus, but when it did appear, it was discussed quite
129 extensively.

130 The term "amendment" is one of the most frequently occurring terms, yet it occurs most
131 frequently in relatively low-weighted topics. And when "amendment" does appear in those
132 topics, it is associated with the terms "fourth," "electronic," "eavesdropping," "privacy,"
133 "surveillance," "evidence," "justice," "telephone," "agent," "conversation," "search,"
134 "seizure," "warrant," and "arrest." While these terms imply discussions of fourth amendment
135 privacy rights vis-à-vis government searches and surveillance, the terms "griswold" and
136 "connecticut" also appear with "amendment," indicating discussion of *Griswold v.*
137 *Connecticut*[1], in which the Supreme Court held that a Connecticut law forbidding the use
138 of contraceptives unconstitutionally intruded upon the right of marital privacy.

139 Based on the 1891-1970 document corpus, three major "areas" of privacy can be discerned
140 in the published literature: rights associated with one's name or likeness, fourth (and more
141 generally ninth and fourteenth) amendment rights against government searches and
142 surveillance, and the emerging issue of information privacy in a rapidly computerizing
143 society.

144

145 **3 Conclusion**

¹ Treemaps were not created for the partial single-year corpora (1980, 1990, 2000, and 2010) due to their limited expanse over time.

146 This paper has presented preliminary results from applying LDA probabilistic topic
147 modeling algorithms to a document collection comprised of published law review articles
148 from 1891 through 1970, and sample publications from 1980, 1990, 2000, and 2010—all
149 citing to Warren’s and Brandeis’s *The Right to Privacy* and substantively discussing privacy
150 law issues. Our initial interpretation of the results reveals a few important trends: early
151 discussion of issues associated with the tort of appropriation of name or likeness; fourth
152 amendment government searches and surveillance; and emerging trends in information
153 privacy associated with data collection, as well as the rise of the Internet. Once our initial
154 data collection has been completed, by filling out the years 1971 through 2012, we hope to
155 not only confirm these initial observations but also reveal additional privacy law trends,
156 developing an ontology of privacy law based upon the topics revealed in published law
157 review and journal articles.

158

159 **Acknowledgments**

160 The authors are extremely grateful for the support and assistance they have received from:
161 Professor Ken Gerow, Chair, University of Wyoming Department of Statistics.

162 **References**

- 163 [1] Griswold v. Connecticut, 381 U.S. 479 (1965).
164 [2] McCallum, A. (2002) *MALLET: A Machine Learning for Language Toolkit* (2002),
165 <http://mallet.cs.umass.edu>.
166 [3] Moore, B., Jr. (1984) *Privacy: Studies in Social and Cultural History*. New York, NY: Pantheon
167 Books.
168 [4] Nissenbaum H. (2010) *Privacy in Context*. Stanford, CA: Stanford Law Books.
169 [5] Roberts, L.A. (1993) *The Ontology of Privacy*. (unpublished Ph.D. dissertation, University of
170 Oregon) (on file with authors).
171 [6] Warren, S.D. & Brandeis, L.D. (1890) The Right to Privacy. *Harvard Law Review* 4(5): 193-220.