

# **Towards Automatic Classification of Privacy Policy Text**

**Frederick Liu    Shomir Wilson    Peter Story  
Sebastian Zimmeck    Norman Sadeh**

June 2018  
CMU-ISR-17-118R  
CMU-LTI-17-010

School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213

This study was supported in part by the NSF under grants CNS-1330596, CNS-1330214, and SBE-1513957. The US Government is authorized to reproduce and distribute reprints for Governmental purposes not withstanding any copyright notation. The views and conclusions contained are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the NSF or the US Government.

**Keywords:** privacy, machine learning, classification, cnn, neural network, privacy policy

## Abstract

Privacy policies notify Internet users about the privacy practices of websites, mobile apps, and other products and services. However, users rarely read them and struggle to understand their contents. Also, the entities that provide these policies are sometimes unmotivated to make them comprehensible. Recently, annotated corpora of privacy policies have been introduced to the research community. They open the door to the development of machine learning and natural language processing techniques to automate the annotation of these documents. In turn, these annotations can be passed on to interfaces (e.g., web browser plugins) that help users quickly identify and understand relevant privacy statements. We present advances in extracting privacy policy paragraphs (termed *segments* in this paper) and individual sentences that relate to expert-identified categories of policy contents, using methods in supervised learning. In particular, we show that relevant segments and sentences can be classified with average micro-F1 scores of 0.78 and 0.66 respectively, improving over prior work. We discuss how the techniques introduced in this paper have been used to automatically annotate the text of about 7,000 privacy policies. Our discussion highlights opportunities as well as limitations associated with our classification approach.



# 1 Introduction

Privacy policies are intended to notify Internet users about the privacy practices that are applicable to their data. Various legal regimes around the world require that website operators, app publishers, and other data processors post a notice on how they gather and share users’ information [12]. This requirement results in a large number of privacy policy documents that most users, however, are unlikely to read. In fact, it was estimated that a user would need to spend at least 181 hours per year to read through all privacy policies for the services they use [5]. It is our goal in this study to make the contents of privacy policies more transparent by leveraging machine learning and natural language processing techniques.

Recent work to annotate large numbers of privacy policies has enabled the use of automated methods toward policy analysis. The OPP-115 Corpus of privacy policies [10], along with its annotation scheme produced by legal experts, provides a springboard for such efforts. Enabled by the OPP-115 Corpus, this paper presents work on using convolutional neural networks (CNNs), logistic regression (LR), and support vector machines (SVMs) to classify policy text into one or more privacy practice *categories*, which represent topics that frequently occur in policy text. We classify policy text at two levels of granularity: sentences and *segments*, which (as defined by the OPP-115 Corpus) roughly correspond to paragraphs. We show CNNs to be competitive on this task with the other two methods, often with higher precision but lower recall. The best results show micro-F1 scores of 0.78 and 0.66 for segment and sentence classification, respectively, suggesting the practicality of tools built upon them. Internet users can benefit in terms of reading less if they are only interested in certain categories. For example, users interested in “First Party Collection/Use” only have to read 16% of the policy at the sentence level and 34% at the segment level as shown in Table 1.

## 2 Related Work

Prior computational work on privacy policy text used information extraction techniques to gather instances of data types mentioned in policies [3], opt-out choices [9], or answers to categorical privacy questions [1, 13, 14]. Closer to our work, one previous study approached the annotation of privacy policy segments as an alignment problem by using Hidden Markov Models [6]. Other

Category	sentence	segment
<b>First Party Collection/Use</b>	373 (16%)	796 (34%)
<b>Third Party Sharing/Collection</b>	280 (12%)	470 (20%)
<b>User, Choice/Control</b>	85 (4%)	164 (7%)
<b>User, Access, Edit &amp; Deletion</b>	34 (2%)	69 (3%)
<b>Data Retention</b>	16 (1%)	15 (1%)
<b>Data Security</b>	54 (2%)	117 (5%)
<b>Policy Change</b>	23 (1%)	67(3%)
<b>Do Not Track</b>	7 (.3%)	12 (1%)
<b>International &amp; Specific Audiences</b>	124 (5%)	194 (8%)

Table 1: Mean absolute counts and percentages of annotated tokens, i.e., words, at the sentence and segment levels per category per policy.

Category	Vocabularies
<b>First Party Collection/Use</b>	use, collect, demographic, address, survey, service
<b>Third Party Sharing/Collection</b>	party, share, sell, disclose, company, advertiser
<b>User Choice/Control</b>	opt, unsubscribe, disable, choose, choice, consent
<b>User Access, Edit and Deletion</b>	delete, profile, correct, account, change, update
<b>Data Retention</b>	retain, store, delete, deletion, database, participate
<b>Data Security</b>	secure, security, seal, safeguard, protect, ensure
<b>Policy Change</b>	change, change privacy, policy time, current, policy agreement
<b>Do Not Track</b>	signal, track, track request, respond, browser, advertising for
<b>International &amp; Specific Audiences</b>	child, California, resident, European, age, parent

Table 2: Vocabulary for each category obtained via logistic regression. Words and collocations are sorted in descending order from left to right according to their weights.

approaches leveraged Latent Dirichlet allocation [2] to facilitate privacy policy comprehension. Our work differs from these prior studies by our framing of the privacy policy analysis task as a multilabel classification problem. This approach is particularly appropriate because a segment of text in a privacy policy can contain information about multiple topics, such as first party collection of data and data security.

The manual annotation of privacy policies has been recognized as a serious bottleneck to modeling their contents, and various prior efforts were aimed at automating the annotation process [7, 11]. They derived motivation from the fact that human annotation is time-consuming, as multiple annotators must carefully interpret legal texts to produce reliable annotations. Some have proposed the automation of assigning category labels to policy segments [10]. Here, we are exploring the OPP-115 Corpus’ use for classifying privacy practices in both sentences and segments.

### 3 OPP-115 Corpus and Annotation Scheme

For our task we make use of the Usable Privacy Policy Project’s [8] OPP-115 Corpus, which contains detailed annotations for the data practices described in a set of 115 website privacy policies [10]. At a high level, annotations fall into one of ten data practice *categories*, which were developed by a team of legal experts:

1. *First Party Collection/Use*: How and why a service provider collects user information
2. *Third Party Sharing/Collection*: How user information may be shared with or collected by third parties
3. *User Choice/Control*: Choices and control options available to users
4. *User Access, Edit, & Deletion*: If and how users can access, edit, or delete their information
5. *Data Retention*: How long user information is stored

6. *Data Security*: How user information is protected
7. *Policy Change*: If and how users will be informed about changes to the privacy policy
8. *Do Not Track*: If and how Do Not Track signals<sup>1</sup> for online tracking and advertising are honored
9. *International & Specific Audiences*: Practices that pertain only to a specific group of users (e.g., children, residents of the European Union, or Californians)
10. *Other*: Additional privacy-related information not covered by the above categories<sup>2</sup>

Privacy policies were divided into *segments*, which were roughly equivalent to paragraphs, and annotators identified spans of text associated with data practices inside of each segment. Each privacy policy was read by three annotators, who required a mean time of 72 minutes per document. In aggregate, they produced a total of 23,194 practice annotations.

We proceed with the observation that the text associated with each category has a distinct vocabulary, even though many of the categories represent closely related themes. Preliminarily, we applied logistic regression to identify particularly relevant words for the different categories. Table 2 shows the results. The top six words or collocations for each category show its distinctiveness.

## 4 Privacy Policy Text Classification

In this section we describe our procedure for labeling privacy policy text at the sentence and segment levels. Different granularity gives different results on the number of tokens annotated, which would result in different reading time if the classification results were used in downstream tasks such as simply highlighting the selected category.

### 4.1 Transforming OPP-115 Annotations into Labels

Annotations for data practices inside a segment can be effectively “elevated” to cover the entire segment, i.e., a segment receives a binary label for the presence or absence of each data practice category. Wilson et al. ([10]) calculated inter-annotator Kappa for segment-level labels to be 0.76 for the first two categories listed above, which comprised 61% of all data practices in the OPP-115 Corpus, and found a variety of lower and higher Kappa values for the remaining categories. For our present work, we use segment-level labels produced by a simple majority vote: if two annotators agree that a segment contains at least one data practice in a given category, then we apply that category to the segment as a label. We use a similar method to produce sentence-level labels: if at least two annotators labeled any part of a sentence with a given category, we label the sentence with that category. Note that the labels are not mutually exclusive, and a segment or sentence may be labeled with zero categories or any combination of them.

---

<sup>1</sup>See [www.w3.org/2011/tracking-protection](http://www.w3.org/2011/tracking-protection).

<sup>2</sup>Because of its indistinct nature, we omit this category from further discussion.

Category	Sentence			Segment			ACL16
	LR	SVM	CNN	LR	SVM	CNN	
First Party Collection/Use	.62/.76/.69	.64/.71/.67	.78/.58/.66	.83/.76/.79	.84/.77/.81	.87/.70/.78	.76/.73/.75
Third Party Sharing/Collection	.57/.73/.64	.61/.72/.66	.86/.40/.55	.71/.85/.77	.74/.81/.78	.79/.80/.79	.67/.73/.70
User Choice/Control	.45/.72/.55	.42/.71/.53	.57/.33/.42	.75/.62/.68	.70/.69/.70	.78/.56/.65	.65/.58/.61
User Access, Edit, & Deletion	.57/.66/.61	.65/.52/.58	.93/.22/.36	.83/.78/.81	.77/.89/.82	.93/.68/.78	.67/.56/.61
Data Retention	.68/.40/.51	.70/.31/.43	.75/.23/.35	.59/.33/.43	.80/.27/.40	0.0/0.0/0.0	.12/.12/.12
Data Security	.62/.74/.67	.60/.71/.65	.67/.71/.69	.67/.79/.73	.70/.85/.77	.77/.85/.80	.66/.66/.67
Policy Change	.66/.80/.72	.75/.78/.77	.86/.65/.74	.95/.74/.83	.95/.67/.78	1.0/.74/.85	.66/.88/.75
Do Not Track	.71/.77/.74	.69/.69/.69	.83/.38/.53	1.0/1.0/1.0	1.0/1.0/1.0	1.0/1.0/1.0	1.0/1.0/1.0
International & Specific Audiences	.75/.74/.74	.75/.75/.75	.77/.69/.73	.72/.86/.79	.88/.82/.85	.79/.84/.81	.70/.70/.70
Average	.61/.73/.66	.63/.70/.66	.78/.51/.60	.77/.76/.76	.80/.77/.78	.80/.71/.75	.66/.66/.66

Table 3: Classification results (precision/recall/F1-score) for sentences and segments using logistic regression (LR), support vector machines (SVM), and convolutional neural networks (CNN). The ACL16 results are from [10]. However, the test set is different from ours and we excluded three classes which belong to the “Other” category under the annotation schema in our average F1 score.

## 4.2 Prediction Methods

For our experiment, we split the 115 policies of the OPP-115 Corpus into 80% training and 20% testing sets. Since each segment or sentence can contain information for multiple categories, we built binary classifiers for each category with three models, respectively logistic regression, support vector machines, and convolutional neural networks [4]. We used a bigram term frequency–inverse document frequency (tf–idf) pre-processor for logistic regression and support vector machines. The parameters for each model are tuned with 5-fold cross validation. The parameters for the CNN follow [4]’s CNN-non-static model, which uses pre-trained word vectors. We used 20% of the training set as a held-out development set to refine these models.

## 4.3 Results and Discussion

The results of segment- and sentence-level classification are shown in Table 3. Across all categories and the best performing model, we observe an average micro-F1 score of 0.78, precision 0.80 and recall 0.77 at the segment level, which outperforms previous results using word-embeddings as features [10].<sup>3</sup> Fewer tokens are annotated at the sentence level. The results at this granularity are on average about 0.1–0.12 worse than at the segment level. One explanation could be that although annotators have access to the context that surrounds a sentence (e.g., prior and subsequent sentences), our sentence-based models do not. We also observe that the CNN model favors precision while the other two models favor recall. This difference can be taken into consideration for downstream tasks with different objectives (e.g., governmental regulators might be interested in manually verifying results; hence, not missing instances is more important than the false positive rate). The results are consistent in F1-scores across the three models as shown in Table 3.

All three models show similar performances after careful parameter tuning, which motivates us to look at the data in more detail to find reasons for errors. For example, the OPP-115 Corpus does not contain many privacy policies of health care providers. One provider’s policy is quoted

<sup>3</sup>The data split may be different since it is not released along with the OPP-115 Corpus.

**A. Third Party Sharing/Collection:**

*As Kaleida Health is a teaching facility, we may disclose your health information for training and educational purposes to faculty physicians, residents and medical, dental, nursing, pharmacy or other students in health-related professions from local colleges or universities affiliated with Kaleida Health.*

**B. Data Security:**

*Chase Paymentech Solutions, LLC is committed to safeguarding the privacy and security of the information we collect.*

Figure 1: Examples of classification errors. For the first example, our models failed to detect the Third Party Sharing/Collection category. In the second example our models disagreed with our gold standard data; however, the text does appear to address Data Security.

in Figure 1A showing health-specific language, more of which would encourage improved performance. We also observed some errors in annotation that may have been oversights by the readers. For instance, the quote in Figure 1B was classified as a security statement, which appears to be correct, but the annotators did not recognize it as such.

During our evaluation we recognized that our classifiers’ performances are also impacted by the context or lack thereof during the production of the annotations. For example, section headings were only shown to the annotators for the segment immediately following, but segments were presented for annotation in order. Features that encode context around each sentence should be investigated to avoid this problem.

Overall, our results indicate the strength of these methods toward enabling downstream tasks, such as filtering for more detailed data practices, extracting salient details to present to users, or summarization of privacy practices.

## 5 Conclusion

In this study we demonstrated the use of traditional and neural network models to classify text in privacy policies according to nine categories that cover important privacy practices. We believe that our results provide support for the use of segment-based annotations. At the same time we recognize that ultimately sentence-based annotations offer the prospect of finer, more detailed annotations. Further research is needed to improve the performance of sentence-based classification. Taking into account the text of adjacent sentences might help. In this paper, we point out the trade-off between number of tokens annotated and classification performance between different granularity of segmenting the policy. For future work, we are developing tools to help reduce the level of effort required from Internet users to understand privacy policies. This includes packaging the results of our analysis in the form of browser plug-ins that summarize key statements extracted from the text of privacy policies, as well as the exploration of question answering functionality to

answer people’s privacy questions.

## Acknowledgments

This study was supported in part by the NSF under grants CNS-1330596, CNS-1330214, and SBE-1513957. The US Government is authorized to reproduce and distribute reprints for Governmental purposes not withstanding any copyright notation. The views and conclusions contained are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the NSF or the US Government.

## References

- [1] Waleed Ammar, Shomir Wilson, Norman Sadeh, and Noah A Smith. Automatic categorization of privacy policies: A pilot study. Technical report, Carnegie Mellon University, 2012. CMU-ISR-12-114, CMU-LTI-12-019.
- [2] Parvathi Chundi and Pranav M. Subramaniam. An approach to analyze web privacy policy documents. In *KDD Workshop on Data Mining for Social Good*, 2014.
- [3] Elisa Costante, Jerry den Hartog, and Milan Petković. What websites know about you: Privacy policy analysis using information extraction. In Roberto Di Pietro, Javier Heranz, Ernesto Damiani, and Radu State, editors, *Data Privacy Management and Autonomous Spontaneous Security*, volume 7731 of *Lecture Notes in Computer Science*, pages 146–159. Springer, 2013.
- [4] Yoon Kim. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*, 2014.
- [5] Aleecia M. McDonald and Lorrie F. Cranor. The cost of reading privacy policies. *I/S: A Journal of Law and Policy for the Information Society*, 4(3):540–565, 2008.
- [6] Rohan Ramanath, Fei Liu, Norman Sadeh, and Noah A. Smith. Unsupervised alignment of privacy policies using hidden markov models. In *Proceedings of the Annual Meeting of the Association of Computational Linguistics*, ACL ’14, pages 605–610. ACL, June 2014.
- [7] Joel R Reidenberg, Travis Breaux, Lorrie Faith Cranor, Brian French, Amanda Grannis, James T Graves, Fei Liu, Aleecia McDonald, Thomas B Norton, and Rohan Ramanath. Disagreeable privacy policies: Mismatches between meaning and users’ understanding. *Berkeley Tech. LJ*, 30:39, 2015.
- [8] Norman Sadeh, Alessandro Acquisti, Travis D Breaux, Lorrie Faith Cranor, Noah A Smith, Fei Liu, Florian Schaub, et al. The usable privacy policy project. Technical report, Carnegie Mellon University, 2013. CMU-ISR-13-119.

- [9] Kanthashree Mysore Sathyendra, Shomir Wilson, Florian Schaub, Sebastian Zimmeck, and Norman Sadeh. Identifying the provision of choices in privacy policy text. In *Conference on Empirical Methods in Natural Language Processing (EMNLP 2017)*, pages 2764–2769, Copenhagen, Denmark, Sep 2017. ACL.
- [10] S Wilson, F Schaub, A Dara, F Liu, S Cherivirala, P G Leon, M S Andersen, S Zimmeck, K Sathyendra, N C Russell, T B Norton, E Hovy, J R Reidenberg, and N Sadeh. The creation and analysis of a website privacy policy corpus. In *Annual Meeting of the Association for Computational Linguistics, Aug 2016*. ACL, 2016.
- [11] Shomir Wilson, Florian Schaub, Rohan Ramanath, Norman Sadeh, Fei Liu, Noah A Smith, and Frederick Liu. Crowdsourcing annotations for websites’ privacy policies: Can it really work? In *Proceedings of the 25th International Conference on World Wide Web*, pages 133–143. International World Wide Web Conferences Steering Committee, 2016.
- [12] Sebastian Zimmeck. The information privacy law of web applications and cloud computing. *Santa Clara Computer & High Tech. L.J.*, 29(3):451–487, 2013.
- [13] Sebastian Zimmeck and Steven M. Bellovin. Privee: An architecture for automatically analyzing web privacy policies. In *23rd USENIX Security Symposium*, USENIX Security ’14, pages 1–16. USENIX Association, August 2014.
- [14] Sebastian Zimmeck, Ziqi Wang, Lieyong Zou, Roger Iyengar, Bin Liu, Florian Schaub, Shormir Wilson, Norman Sadeh, Steven M. Bellovin, and Joel Reidenberg. Automated analysis of privacy requirements for mobile apps. In *24th Network & Distributed System Security Symposium (NDSS 2017)*, NDSS 2017, San Diego, CA, February 2017. Internet Society.