The University of Texas at Austin
# Center for Identity

# Privacy: a Machine Learning Perspective

*Razieh Nokhbeh Zaeem*
*K. Suzanne Barber*

# Economics of Privacy: Privacy, a Machine Learning Perspective

Razieh Nokhbeh Zaeem* and K. Suzanne Barber

## Synonyms

Machine Learning (Data Mining) for Privacy Policy Summarization.

## Definitions

Machine Learning for Privacy Policy Summarization is the novel application of machine learning techniques to automatically extract summaries of online privacy policies.

## Background

Privacy policies have become the de facto way of communicating how a company or organization–and particu-

larly its website–collects, shares, and uses personally identifiable information (PII). These privacy policies outline how the organization handles, shares, discloses, and uses PII of its consumers or clients. PII is defined as "any information relating to an identified or identifiable natural person"[2] such as name, email address, and credit card number.

Meinert et al (2006) showed that while most users know about privacy policies, less than half of them have *ever* read a privacy policy. Milne and Culnan (2004) used self-reported data from users and found that only 4.5% claim to always read them. However, using the more reliable server side observation of websites, Kohavi (2001) revealed even more astonishing statistics that only 1% or less of users click on a website's privacy policy. More recently, Steinfeld (2016) used advanced eye tracking tech-

* corresponding author

[2] https://gdpr-info.eu

niques to demonstrate that the same still holds true today: users barely take effort to read privacy policies thoroughly.

Ermakova et al (2014) attribute the fact that most users do not read privacy policies to the privacy policies' poor readability. Graber et al (2002) and Milne et al (2006) show that the average privacy policy requires two years of college level education to comprehend. In addition, Milne et al (2006) found that privacy policies are getting longer and harder to read, with the readability score of the privacy policies decreasing over time. In fact, according to McDonald and Cranor (2008), reading privacy policies is so time consuming that if users were to read each new privacy policy they encounter in a year, it would take them over 200 hours.

To address the lack of readability in privacy policies, researchers have developed tools that leverage Machine Learning (ML) and data mining to automatically summarize conventional privacy policies for the users. Among these tools, however, very few are made available to the public.

## Application

Privee by Zimmeck and Bellovin (2014) is the first automatic privacy policy analysis tool. Building on the crowd sourcing privacy analysis framework ToS;DR (2012), Privee combines crowd sourcing with rule and machine learning classifiers to classify privacy policies that are not already rated in the crowd sourcing repository.

The Usable Privacy Project by Sadeh et al (2013) takes advantage of natural language processing, machine learning, privacy preference modeling, crowd sourcing, and formal methods to semi-automatically annotate privacy policies.

Two of the most recent publicly available tools to utilize machine learning for privacy policy summarization are Polisis[3] by Harkous et al (2018) and PrivacyCheck[4] by Zaeem et al (2018).

At its core, Polisis is a neural network classifier trained on 130,000 privacy policies retrieved from the Google Play store. Polisis segments a privacy policy and automatically annotates each segment with a set of labels, classifying segments based on coarse- and fine-grained classifications. Pribots by Harkous et al (2016) is from the same authors and is a chat bot that answers free-form questions about privacy policies.

Powered by machine learning, PrivacyCheck by Nokhbeh Zaeem and Barber (2020) is a publicly available browser extension that automatically summarizes any privacy policy by answering 20 questions based upon User Control and the General Data Protection Regulation (GDPR). Furthermore, PrivacyCheck incorporates a competitor analysis tool that highlights the top competitors with the best privacy policies in the same market sector. PrivacyCheck enhances the users' understanding of privacy policies and empowers them to make informed decisions when it comes to selecting services with better privacy policies.

Other researchers, too, have applied machine learning and natural language processing in privacy policy analy-

[3] Available online at https://pribot.org/polisis.
[4] Available online at https://identity.utexas.edu/privacycheck-for-google-chrome.

sis Clarke et al (2012); Fawaz et al (2019). PolicyLint by Andow et al (2019) is a natural language processing tool that identifies potential contradictions that may arise inside the same privacy policy. PolicyLint is tested on a corpus of 11,430 privacy policies from mobile apps. PrivacyGuide by Tesfay et al (2018) is a machine learning and natural language processing tool inspired by the GDPR. It uses a corpus of 45 policies from the most accessed websites in Europe. PrivacyGuide, however, is not publicly available, as opposed to PrivacyCheck and Polisis. Finally, studies of mobile app privacy policies (e.g., MAPS by Zimmeck et al (2019)) are on the rise.

### *Economics of Privacy*

Acquisti et al (2016) reviewed the economics of privacy and reported that, in digital economies, consumers often receive imperfect, incorrect, or asymmetric information regarding what data is collected about them and how that data will be used. Hence, consumers' ability to make informed decisions about their privacy is severely hindered. Data mining tools that summarize privacy policies directly address this need and seek to improve users' understanding of what they agree to in a privacy policy.

## Open problems and Future directions

A prominent future work is to gather a more comprehensive training dataset, in order to train more accurate supervised machine learning models. Finally, encouraging widespread use of such tools by final consumers is an important future direction that should be pursued.

## Cross-References

1. Privacy-preserving data mining: Data Mining (Privacy in)

2. Privacy metrics and data protection: Personally Identifiable Information

3. Privacy metrics and data protection: Privacy-Enhancing Technologies

## References

Acquisti A, Taylor C, Wagman L (2016) The economics of privacy. Journal of Economic Literature 54(2):442–492

Andow B, Mahmud SY, Wang W, Whitaker J, Enck W, Reaves B, Singh K, Xie T (2019) Policylint: investigating internal privacy policy contradictions on Google play. In: 28th USENIX Security Symposium (USENIX Security 19), pp 585–602

Clarke N, Furnell S, Angulo J, Fischer-Hübner S, Wästlund E, Pulls T (2012) Towards usable privacy policy display and management. Information Management & Computer Security 20(1):4–17

Ermakova T, Baumann A, Fabian B, Krasnova H (2014) Privacy policies and users' trust: Does readability matter? In: 20th Americas Conference on Information Systems (AMCIS)

Fawaz K, Linden T, Harkous H (2019) The applications of machine learning in privacy notice and choice. In: 2019 11th International Conference on Communication Systems & Networks (COMSNETS), IEEE, pp 118–124

Graber MA, D Alessandro DM, Johnson-West J (2002) Reading level of privacy policies on internet health web sites. Journal of Family Practice 51(7):642–642

Harkous H, Fawaz K, Shin KG, Aberer K (2016) Pribots: Conversational privacy with chatbots. In: Twelfth Symposium on Usable Privacy and Security (SOUPS) 2016

Harkous H, Fawaz K, Lebret R, Schaub F, Shin KG, Aberer K (2018) Polisis: Automated analysis and presentation of privacy policies using deep learning. In: 27th USENIX Security Symposium (USENIX Security 18), pp 531–548

Kohavi R (2001) Mining e-commerce data: the good, the bad, and the ugly. In: International conference on Knowledge discovery and data mining, ACM, pp 8–13

McDonald AM, Cranor LF (2008) the cost of reading privacy policies. I/S: A Journal of Law and Policy for the Information Society 4:543

Meinert DB, Peterson DK, Criswell JR, Crossland MD (2006) Privacy policy statements and consumer willingness to provide personal information. Journal of Electronic Commerce in Organizations 4(1):1

Milne GR, Culnan MJ (2004) Strategies for reducing online privacy risks: Why consumers read (or don't read) online privacy notices. Journal of Interactive Marketing 18(3):15–29

Milne GR, Culnan MJ, Greene H (2006) A longitudinal assessment of online privacy notice readability. Journal of Public Policy & Marketing 25(2):238–249

Nokhbeh Zaeem R, Barber KS (2020) The effect of the gdpr on privacy policies: Recent progress and future promise. ACM Transactions on Management Information Systems (TMIS)

Sadeh N, Acquisti A, Breaux TD, Cranor LF, McDonalda AM, Reidenbergb JR, Smith NA, Liu F, Russellb NC, Schaub F, et al (2013) The usable privacy policy project. Tech. rep., Technical Report, CMU-ISR-13-119, Carnegie Mellon University

Steinfeld N (2016) "I agree to the terms and conditions": (how) do users read privacy policies online? an eye-tracking experiment. Computers in Human Behavior 55:992–1000

Tesfay WB, Hofmann P, Nakamura T, Kiyomoto S, Serna J (2018) Privacyguide: towards an implementation of the EU GDPR on internet privacy policy evaluation. In: Proceedings of the Fourth ACM International Workshop on Security and Privacy Analytics, ACM, pp 15–21

ToS;DR (2012) Terms of service; didn't read. URL https://tosdr.org

Zaeem RN, German RL, Barber KS (2018) Privacycheck: Automatic summarization of privacy policies using data mining. ACM Transactions on Internet Technology (TOIT) 18(4):53

Zimmeck S, Bellovin SM (2014) Privee: An architecture for automatically analyzing web privacy policies. In: 23rd USENIX Security Symposium (USENIX Security 14), USENIX Association, San Diego, CA, pp 1–16

Zimmeck S, Story P, Smullen D, Ravichander A, Wang Z, Reidenberg J, Russell NC, Sadeh N (2019) Maps: Scaling privacy compliance analysis to a million apps. Proceedings on Privacy Enhancing Technologies 2019(3):66–86