

contributed articles

DOI:10.1145/2699410

Legitimacy of surveillance is crucial to safeguarding validity of OSINT data as a tool for law-enforcement agencies.

BY PETRA SASKIA BAYERL AND BABAK AKHGOR

Surveillance and Falsification Implications for Open Source Intelligence Investigations

OPEN SOURCE INTELLIGENCE, or OSINT, has become a permanent fixture in the private sector for assessing consumers' product perceptions, tracking public opinions, and measuring customer loyalty.¹² The public sector, and, here, particularly law-enforcement agencies, including police, also increasingly acknowledge the value of OSINT techniques for enhancing their investigative capabilities and response to criminal threats.⁵

OSINT refers to the collection of intelligence from information sources freely available in the public

domain, including offline sources (such as newspapers, magazines, radio, and television), along with information on the Internet.^{4,16,17} The spread of social media has vastly increased the quantity and accessibility of OSINT sources.^{3,11} OSINT thus complements traditional methods of intelligence gathering at very low to no cost.^{4,15}

OSINT increasingly supports the work of law-enforcement agencies in identifying criminals and their activities (such as recruitment, transfer of information and money, and coordination of illicit activities);¹⁸ for instance, the capture of Vito Roberto Palazzolo, a treasurer for the Italian mafia on the run for 30 years was accomplished in part by monitoring his Facebook account.⁸ OSINT also demonstrated its potential to help respond quickly to criminal behavior outside the Internet, as during, for instance, public disorder (such as the 2011 U.K. riots).¹ OSINT has therefore become an important tool for law-enforcement agencies combating crime and ultimately safeguarding our societies.¹⁴

To fulfill these functions, OSINT depends on the integrity and accuracy of open data sources. This integrity is jeopardized if Internet users choose not to disclose personal information or even provide false information on

» key insights

- **Falsification of personal information online is widespread, though users do not falsify information in a uniform way; some types of information are more likely to be falsified, leading to systematic differences in the reliability of the information for OSINT applications.**
- **Acceptance of and propensity for falsification is linked to attitudes toward online surveillance; the more negative people's attitudes toward governmental online surveillance, the more likely they are to accept information falsification and provide false information.**
- **The source of online surveillance—state agencies vs. private companies vs. unnamed organizations—seems to affect whether assumptions about online surveillance are linked to information falsification.**



themselves.^{7,9} Such omissions and falsifications can have grave consequences if decisions are being made from data assumed to be accurate but that is not.¹⁹

This issue is especially relevant since the revelations by former NSA contractor Edward Snowden of large-scale monitoring of communications and online data by state agencies

worldwide. The revelations created considerable mistrust by citizens of Internet-based surveillance by their own governments, bringing the tension between the security of society and the fundamental right to privacy into sharp profile. These discussions begin to show concrete effects; for instance, use of privacy-sensitive keywords in Google searches changed

from the period before to the period after Snowden's revelations, as users proved less willing to use keywords "that might get them in trouble with the [U.S.] government."¹⁰ Despite mandatory national and international data protection and privacy regulations, Internet users thus seem wary of online surveillance and in consequence modify their online behavior.

For organizations using OSINT in their decision making, changes in users' online behaviors, specifically their willingness to provide accurate accounts about themselves, are problematic. Not only do they increase the incidence of false information, they also raise the complexity and costs for information validation, or authentication of individuals' Web footprints, against additional and trusted sources.

A better understanding of the tendency of Internet users of when and why to change their online behavior in response to online surveillance can help pinpoint especially problematic areas for the validity of OSINT methods. Such an understanding can further guide efforts for more targeted cross validations. So far, organizations, including law-enforcement agencies, lack a clear

picture of how far and in what ways concerns about online surveillance change information bases relevant to law-enforcement agencies' use of open source intelligence. In our current research, we aim to systematically investigate whether shifts in online behaviors are likely and, if so, what form they might take. In this article, we report on a study in which we focused on the falsification of personal information, investigating the link between falsification acceptance and falsification propensity with attitudes toward online surveillance, privacy concerns, and assumptions about online surveillance by different organizations.

Study Design and Sample

To understand Internet users' attitudes toward falsification of personal information in connection with online surveillance, we conducted an online survey between January and March 2014 using the micro-work platform Amazon Mechanical Turk to recruit participants.^a A total of 304 users responded to our request, of which 298 provided usable answers. Our sample consisted largely of experienced Internet users (72.2% had more than 11 years of experience) and intensive users (41.3% using the Internet for at least seven hours per day). The majority (83.9%) of participants lived in the U.S., 9.4% in India, and the others were from Canada, Croatia, Kenya, or Romania (0.4% to 1.1% per country). The gender distribution was nearly equal, with 48.9% male vs. 50.4% female participants; 0.7% preferred not to answer the question. Participants were relatively young, with a majority 40 years or younger (67.3%), of which most (35.6%) were between 21 and 30 years of age. Older participants were slightly underrepresented, with 9.5% between 51 and 60 and 3.9% over 60; 0.7% preferred not to answer the question. The questionnaire was administered online. On completion of the survey participants were paid \$0.70 through the Mechanical Turk platform. The survey took an average four minutes to complete.

^a Amazon Mechanical Turk is an online service that allows recruitment of participants worldwide for jobs of short duration, often lasting only several minutes; see <https://www.mturk.com>

Figure 1. Respondents' attitudes toward the positive and negative sides of state online surveillance.

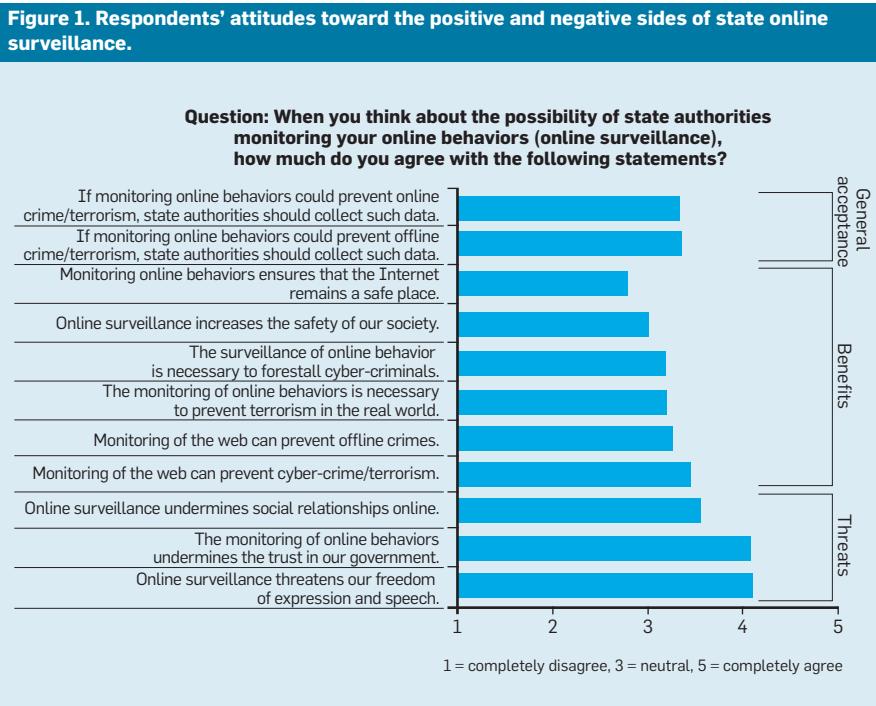
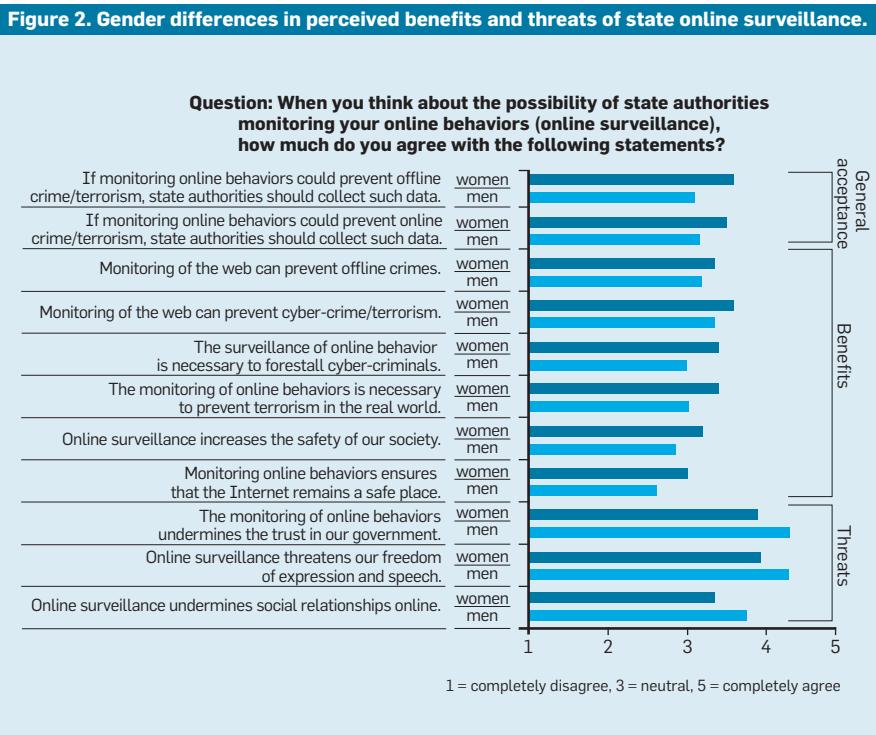


Figure 2. Gender differences in perceived benefits and threats of state online surveillance.



Findings

In the following sections, we detail our findings on participants' attitudes toward surveillance, acceptance of and propensity for falsification of their personal online information, and the possible links between them.

Attitudes toward online surveillance by state agencies. The first question when investigating the effect of state surveillance on online behaviors is how Internet users perceive its value. To capture attitudes toward online surveillance by state agencies we asked participants to indicate their agreement with 11 statements, five of them positive toward online surveillance, thus addressing potential benefits, three of them negative, addressing possible threats, and two capturing general acceptance; Figure 1 shows the average values for benefits, threats, and general acceptance for the entire sample.

The general acceptance of online surveillance was at a medium level with $m = 3.35$ when the focus was on the prevention of offline crimes and $m = 3.33$ when focusing on the prevention of online crimes (both measured on a scale of 1 to 5). Overall, negative attitudes were considerably stronger than positive attitudes. Participants were especially concerned about threats to freedom of expression and speech and the undermining of trust in their own government. Interestingly, the claims state agencies often make that monitoring online behavior ensures the Internet stays a safe place or increases the safety of society found little agreement.

Women were generally more accepting of online surveillance ($t(280) = -3.02$, $p <.01$), seeing significantly more benefits than men ($t(279) = -2.60$, $p <.01$). Men in contrast reported significantly more concern about its negative aspects ($t(275) = 3.69$, $p <.001$) (see Figure 2). Women were especially more willing to support online surveillance if it could prevent crimes perpetrated outside the Internet (offline crimes), whereas men were particularly concerned about the undermining of trust in the government. Moreover, users with more experience in the use of the Internet (more than 11 years) were significantly less positive toward online surveillance than users with less ex-

Figure 3. Respondents' assumptions concerning the degree of online surveillance by different organizations.

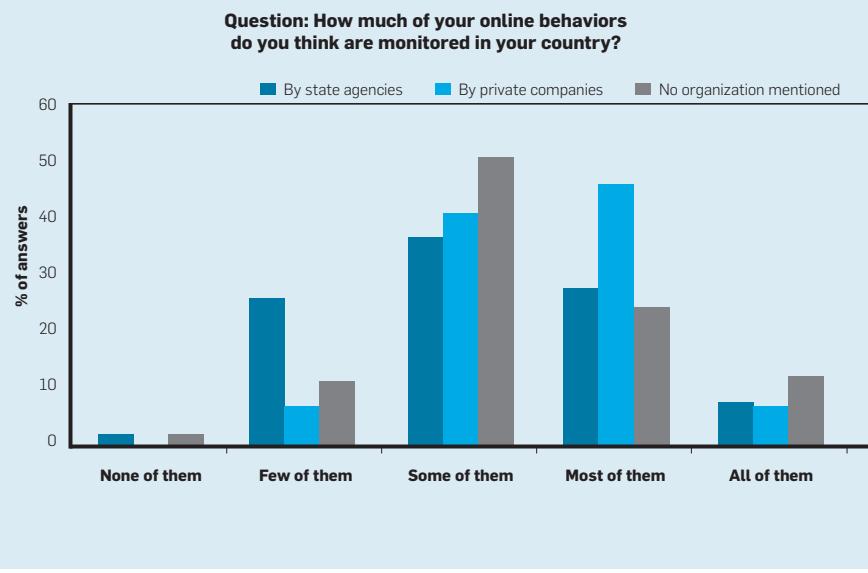
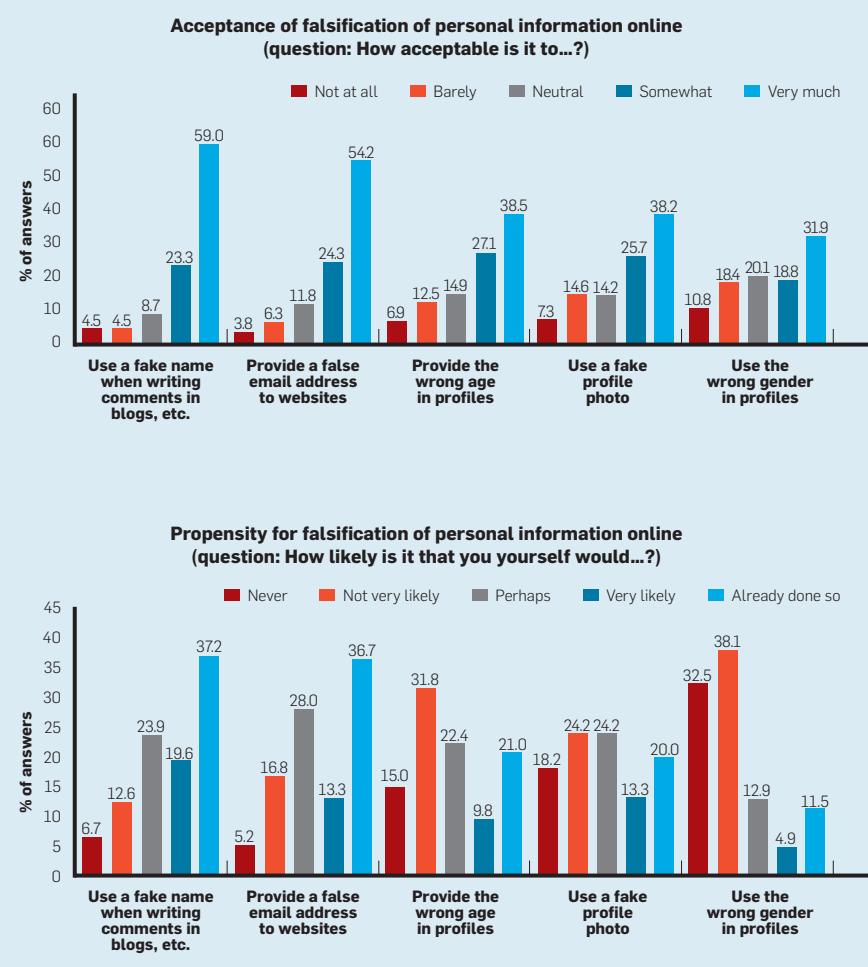


Figure 4. Acceptance and propensity for falsification of personal information among all participants.



perience (seven years or less; $F(2,274) = 5.04, p <.01$). Since age groups did not differ in their attitudes, this effect cannot be explained by generational differences. It instead hints at growing sensitivity toward the issue with increased Internet use.

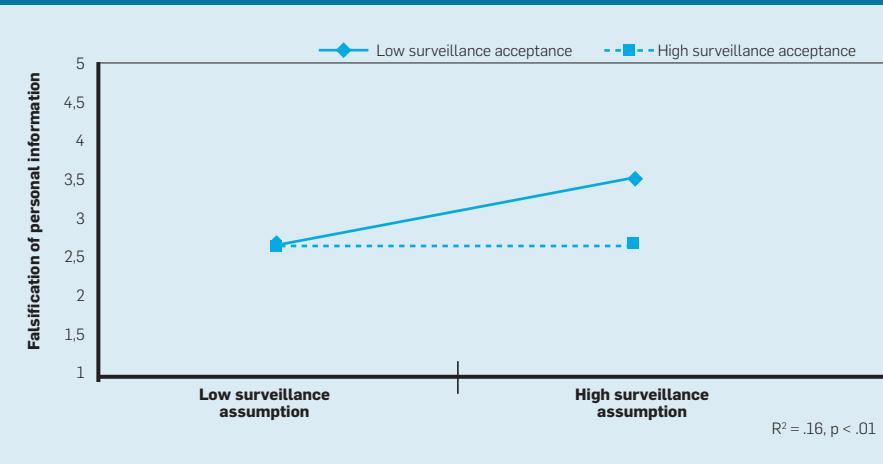
Surveillance by state agencies vs. private companies or unnamed organizations. Unlike private companies, which are widely known for collecting online data on a large scale, OSINT use by state agencies has only recently come to the attention of the broader public. Yet, as demonstrated by the intense discussion in the aftermath of the Snowden

revelations, the sensitivity of the issue seems even greater. Also, compared to the use of OSINT by private companies, the consequences of OSINT use by law-enforcement agencies can be considerably more severe for the individual under scrutiny. We therefore wanted to know whether online surveillance by state agencies could lead to different reactions from surveillance conducted by private industry. For the second part of the survey we used three different framings for our questions, one mentioning surveillance conducted by state agencies, one mentioning surveillance by private companies, and one men-

tioning surveillance without naming a specific organization. A total of 104 people, or 34.9% of respondents, filled out the survey referring to surveillance by state authorities; 103 people (34.6%) answered the survey referring to surveillance by private companies; and 91, or 30.5%, reacted to the generalized condition in which no specific organization was named.

First, we were interested in the extent of online surveillance users assumed across the three sources of surveillance, ranging from “none” of their online behaviors to “all of them.” In all three conditions, the average values indicate users assumed at least some of their behavior is monitored, although the values were highest for private companies ($m = 3.52$) and lowest for state agencies ($m = 3.13$) (see Figure 3). This difference was also statistically significant ($F(2,294) = 5.37, p <.01$). This was a general tendency, as genders, age groups, and user groups with different degrees of Internet experience did not differ in their assumptions about online surveillance. Despite current debates, private companies thus seem to be perceived as more intrusive than state agencies. As we outline in the following sections, this does not mean, however, that surveillance by state

Figure 5. Role of surveillance assumptions and acceptance in information falsification.



Correlations between falsification behaviors and online surveillance assumptions and attitudes.

Generic condition (no mention of an organization; n = 91)

	Mean	Std. dev.	1.	2.
1. Assumption of online surveillance	3.36	0.88		
2. Acceptance of information falsification	3.80	1.06	.22	
3. Propensity for information falsification	3.02	1.03	.10	.66**

Condition “surveillance by private companies” (n = 103)

	Mean	Std. dev.	1.	2.
1. Assumption of online surveillance	3.52	0.73		
2. Acceptance of information falsification	3.99	0.96	.13	
3. Propensity for information falsification	3.26	1.03	.12	.63**

Condition “surveillance by state agencies” (n = 104)

	Mean	Std. dev.	1.	2.	3.	4.	5.
1. Assumption of online surveillance	3.13	0.96					
2. General acceptance of online surveillance by state agencies	3.23	1.22	-.04				
3. Benefits from surveillance	3.06	1.02	.01	.78**			
4. Threats from surveillance	4.05	0.79	.11	-.38**	-.49**		
5. Acceptance of information falsification	3.84	0.96	.08	-.32**	-.24**	.21*	
6. Propensity for information falsification	2.92	1.07	.24*	-.26**	-.23*	.13	.59**

* $p < .05$

** $p < .01$; Pearson correlations, two-sided tests

agencies is seen as less intrusive than that of private companies.

Degree of acceptance and propensity to falsify personal information online.

To understand whether concerns about online surveillance affect the tendency to falsify personal information online, we asked participants in all three conditions the same two questions: How acceptable is it to falsify personal online information (acceptance of falsification, from 1 = not at all to 5 = very much)? And how likely are you to falsify your own personal information online (propensity for falsification, from 1 = never to 5 = already done so)?

We asked participants about the falsification of five types of information that are fixtures in most online profiles: providing a false name, providing a fake email address, providing the wrong age, using a fake photo, and providing the wrong gender.

Taking all five together, users showed a high level of acceptance for falsification ($m = 3.88$, $SD = 0.99$), while the propensity for falsification was somewhat less ($m = 3.06$, $SD = 1.05$). Still, only a very small group (3.4%) indicated they would never falsify any information, whereas 7.4% indicated having already done so for all five categories.

Interestingly, falsification acceptance and propensity were not uniform across all five. Using a false name and false email address was seen as acceptable, whereas a false profile photo and wrong gender were considered much less acceptable (see top part of Figure 4). Only 9.0% of the participants considered falsifying their own name as completely or highly unacceptable; for the falsification of one's gender, this was 29.2%. The same trend emerged for the propensity of falsifying information; 37.0% of participants indicated they had already used a fake name and email address, while 70.6% reported they would never use the wrong gender or were very unlikely to do so (see bottom part of Figure 4). Users thus seem nearly five times more likely to indicate the wrong name and more than six times more likely to provide a wrong email address than report the wrong gender. This suggests the falsification of personal information follows specific patterns; that is, dif-

The more participants perceived online surveillance by state agencies as problematic, the more willing they were to accept falsification.

ferent pieces of information in a profile may have a disparate likelihood of being valid or invalid, or "differential validity of information types."

To compare the effect of the three surveillance sources, we summarized the five types of information into one score for acceptance and one score for propensity, respectively. The three conditions did not differ in terms of falsification acceptance ($F(2,285) = 0.92$, nonsignificant) but resulted in at least a marginal effect for falsification propensity ($F(2,281) = 2.77$, $p = .06$). This was due to a slightly greater propensity for falsification when surveillance was conducted by private companies ($m = 3.26$) compared to state agencies ($m = 2.91$; $t = -2.29$, $p < .05$). Gender, age groups, and length of Internet use had no effect on either outcome.

Linking information falsification with surveillance assumptions and attitudes. We next considered influence of surveillance awareness, attitudes toward surveillance, and privacy concerns on information falsification. Because we used three separate versions of the survey to determine the influence of the organization conducting surveillance, the questions on degree of surveillance awareness and falsification acceptance and propensity referred to different entities: state agencies, private organizations, or no organization in particular. We therefore calculated the correlations between surveillance awareness and information falsification for each of the three groups separately. This also gave us the opportunity to investigate whether the context of surveillance had an effect on falsification behaviors. The table here reports the results for each of the three conditions.

Interestingly, assumptions of online surveillance had an effect on falsification acceptance and propensity only when online surveillance was framed in the context of state agencies or as generalized activity. In these cases, assumptions about online surveillance had a clear positive link with either the propensity to falsify personal information or the acceptance of this behavior, as in the table. For surveillance conducted by private companies, no such significant link emerged. Again, this suggests the question of who conducts the sur-

veillance may play a role in influencing concrete falsification behaviors. Surveillance by state agencies could trigger more concrete reactions than either generalized surveillance or monitoring by private companies.

As in the third generalized condition, all questions referred uniformly to state agencies, this subgroup of participants gave us the opportunity to further investigate the link between attitudes toward online surveillance by state agencies and falsification. In this subgroup, we found a clear link between attitudes toward online surveillance, acceptance, and propensity for falsification. The greater the general acceptance and perceived benefits of surveillance, the less accepting participants were of falsifying information and the less likely they were to do it themselves. Similarly, the more participants perceived online surveillance by state agencies as problematic, the more willing they were to accept falsification.

In addition, acceptance of online surveillance moderated the relationship between falsification and assumed degree of surveillance. While greater assumptions of surveillance generally increased the propensity for falsification, this reaction was especially strong for people with a low acceptance of online surveillance by state agencies (see Figure 5). This observation suggests an important interaction between awareness and attitudes. While surveillance awareness alone may lead to information falsification, the main trigger to falsifying personal information seems to be the extent surveillance is seen as (in)appropriate. This logic links tendencies for falsification of one's own information to how much one considers state agencies legitimate and trustworthy, thus emphasizing the potentially critical effect of negative press on the viability of OSINT-based decisions.

More than a Moral Dilemma

Our study demonstrates that discussions about "privacy" vs. the "rightfulness" of online surveillance is more than a moral dilemma. Rather, the degree to which individuals are aware of online surveillance and the way they view the acceptability of this act, including the organizations implicated

Law-enforcement agencies will have to become more sensitive to the reactions their own practices might create for the viability of their methods and in consequence the decisions they take based on these methods.

in it, can pose concrete challenges for the validity of online data—and consequently for the validity of decisions based on the data. While our study is only a small window into this complex issue, it demonstrates that online surveillance may have very concrete, practical implications for the use and usefulness of OSINT, specifically for law-enforcement agencies. Surveillance is not neutral. On the contrary, our study attests that surveillance practices could threaten the integrity of the very data they rely on.

Falsification tendencies as a reaction to online surveillance create challenges for the usability of open source data, especially increasingly for the effort required to validate information. OSINT has long been hailed as a cheap or even "no cost" source of operational information for law-enforcement agencies.^{4,16} Our findings suggest that increasing awareness of online surveillance, including painful revelations of problematic surveillance practices by states and law-enforcement agencies, may severely reduce this benefit, at least for those Internet users with a more critical outlook toward state authorities and/or greater need for privacy.

Technical solutions to counter the increased likelihood of falsification are available; for instance, Dai et al.⁵ proposed a number of "trust score" computation models that try to determine data trustworthiness in anonymized social networks using a trusted standard. Additional solutions are thinkable using validity pattern mining, reasoning-based semantic data mining, and open source analysis techniques. One important avenue for identifying false information is to identify possible links between profiles of a single user and then mine the data between profiles for validation. Users often explicitly link their profiles. For example, Twitter posts and Instagram photos can be organized so they appear on a user's Facebook timeline. This gives a direct and verified link to further information. Users may also post under the same pseudonym on a number of profiles. Collecting the data associated with each of these profiles provides further opportunity for corroboration.

As with Dai et al.,⁵ another tactic might be to attempt to match the social graph of users across networks. Inconsistencies in personal data may be identified by verifying where these networks overlap.

The most difficult part in information validations is determining the technological solutions that must be employed to carry out the validation. Two such techniques are classification and association mining. Machine-learning-based classification techniques can be used to establish a ground-truth dataset containing information known to be accurate. By training models on this data, outliers in new data could indicate the trustworthiness of the information may warrant further investigation. Association mining (or association rule learning) can be used to discover relationships between variables within datasets, including social media and other OSINT sources.¹² These association rules can take data from the links discovered between multiple social networks and be used to validate the existing information.

Still, all these technical solutions rely on the cross-validation of open source information with other (open or closed) sources. Growing falsification tendencies in the wake of increasing online surveillance awareness will make such cross-validations not only increasingly necessary but also more complex and costly. Here, the notion of differential validity, as evidenced in our data, may provide a valuable perspective toward a more systematic and targeted approach to information validation by guiding validation efforts toward more or less problematic data. This approach follows the observation that personal information seems to possess systematic variations in its veracity, leading to differential validity patterns. While our study focused on only a small set of static personal information, we assume similar patterns are also observable for other areas, as well as for more dynamic data.

An interesting question in this regard is how “volatile” falsifications of personal information tend to be. Do users stick with one type of falsification (such as consistently modifying

name, relationship status, or age) across services, or do these pieces of information vary across services? Also, do users always use the same content (such as the same false date of birth or photo)? Extending our knowledge of such falsification or validity patterns can considerably reduce the effort involved in validating OSINT-based data. In our current study, we did not investigate the reasons behind the differences in falsification acceptance and propensity for the various types of personal information. Getting a clearer understanding of these reasons could tell us much about the contexts in which falsification are more or less likely, as well as the strategies Internet users employ to remain private.

Conclusion

We clearly cannot return to the days of the “uninformed” or “unaware” Internet user, and law-enforcement agencies therefore need to find ways to deal with the consequences of online surveillance awareness by the general public and the possible ramifications it may have for the trustworthiness of online information. While we do not suggest OSINT will lose its value for investigation processes, we certainly think law-enforcement agencies will have to become more sensitive to the reactions their own practices might create for the viability of their methods and in consequence the decisions they take based on these methods.

Employing ever more advanced technical solutions is not the (sole) solution. Our findings make clear that even more than the pure fact of online surveillance, it is the perceived purpose and legitimacy of the act that are the main drivers behind the extent to which users alter their behaviors online. This explains the role of (largely negatively tinted) public discussions for the behavioral changes in the wake of Snowden’s revelations.¹⁰ They also highlight the criticality of properly legitimizing online surveillance to reduce distrust in law-enforcement agencies and thus pressures toward information falsifications and probably changes in online behaviors more generally.

References

- Barlett, J., Miller, C., Crump, J., and Middleton, L. *Policing in an Information Age*. Demos, London, U.K., Mar. 2013.
- Bell, P. and Congram, M. Intelligence-led policing (ILP) as a strategic planning resource in the fight against transnational organized crime (TOC). *International Journal of Business and Commerce* 2, 12 (2013), 15–28.
- Best, C. Challenges in open source intelligence. In *Proceedings of the Intelligence and Security Informatics Conference* (Athens, Greece, Sept. 12–14, 2011), 58–62.
- Best Jr., R.A. and Cumming, A. *Open Source Intelligence (OSINT): Issues for Congress*. Congressional Research Service, Washington, D.C., Dec. 2007; <https://www.fas.org/sgp/crs/intel/RL34270.pdf>
- Dai, C., Rao, F.Y., Truta, T.M., and Bertino, E. Privacy-preserving assessment of social network data trustworthiness. In *Proceedings of the Eighth International Conference on Collaborative Computing* (Pittsburgh, PA, Oct. 14–17, 2012), 97–106.
- Gibson, S. Open source intelligence: An intelligence lifeline. *The RUSI Journal* 149, 1 (2004), 16–22.
- Joinson, A.N., Reips, U.D., Buchanan, T., and Schofield, C.B.P. Privacy, trust, and self-disclosure online. *Human-Computer Interaction* 25, 1 (2010), 1–24.
- La Stampa. Mafia, fermato Vito Roberto Palazzolo scovato a Bangkok grazie a Facebook. *La Stampa* (Mar. 31, 2012); <http://www.lastampa.it/2012/03/31/italia/cronache/mafia-fermato-vito-roberto-palazzoloscovato-a-bangkok-grazie-a-facebook-vpxnhM5z5chH3iu1jttsJ/pagina.html>
- Lenhart, A., Madden, M., Cortesi, S., Duggan, M., Smith, A., and Beaton, M. *Teens, Social Media, and Privacy*. Pew Internet and American Life Project Report, Washington, D.C., 2013; <http://www.pewinternet.org/2013/05/21/teens-social-media-and-privacy/>
- Marthew, A. and Tucker, C. *Government Surveillance and Internet Search Behavior*. Working Paper. Social Science Research Network, Rochester, NY, Mar. 2014; https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2412564
- Mercado, S.C. Sailing the sea of OSINT in the information age. *Studies in Intelligence* 48, 3 (2009), 45–55.
- Nancy, P., Ramani, R.G., and Jacob, S.G. Mining of association patterns in social network data (Facebook 100 universities) through data mining techniques and methods. In *Proceedings of the Second International Conference on Advances in Computing and Information Technology*. Springer, Berlin, 2013, 107–117.
- Neri, F., Aliprandi, C., Capaci, F., Cuadros, M., and By, T. Sentiment analysis on social media. In *Proceedings of the 2012 International Conference on Advances in Social Networks Analysis and Mining* (Istanbul, Turkey, Aug. 26–29, 2012), 919–926.
- Omand, D., Bartlett, J., and Miller, C. Introducing social media intelligence. *Intelligence and National Security* 27, 6 (2012), 801–823.
- Ratzel, M.P. Europol in the combat of international terrorism. *NATO Security Through Science Series, Volume 19*. IOS Press, Amsterdam, 2007, 11–16.
- Steele, R.D. The importance of open source intelligence to the military. *International Journal of Intelligence and Counter Intelligence* 8, 4 (1995), 457–470.
- Steele, R.D. Open source intelligence. Chapter 10 in *Handbook of Intelligence Studies*, J. Loch, Ed. Routledge, New York, 2007, 129–147.
- Stohl, M. Cyberterrorism: A clear and present danger, the sum of all fears, breaking point, or patriot games? *Crime, Law, and Social Change* 46, 4–5 (2006), 223–238.
- The Telegraph. Connecticut school shooting: Police warn of social media ‘misinformation.’ *The Telegraph* (Dec. 16, 2012); <http://www.telegraph.co.uk/telegraphtv/9748745/Connecticut-school-shooting-police-warn-of-social-media-misinformation.html>

Petra Saskia Bayerl (pbayerl@rsm.nl) is an assistant professor for technology and organizational behavior and program director technology of the Center of Excellence in Public Safety Management at the Rotterdam School of Management, Erasmus University Rotterdam, the Netherlands.

Babak Akhgar (B.Akhgar@shu.ac.uk) is a professor of informatics and director of the Center of Excellence in Terrorism, Resilience, Intelligence, and Organized Crime Research at Sheffield Hallam University, UK.