# New tool expands tracking of personal data on the Web

October 9 2015



A second-generation tool called Sunlight is intended to bring greater transparency to the Web. Credit: (Tumblr)

Navigating the Web gets easier by the day as corporate monitoring of

our emails and browsing habits fine-tune the algorithms that serve us personalized ads and recommendations. But convenience comes at a cost. In the wrong hands, our personal information can be used against us, to discriminate on housing and health insurance, and overcharge on goods and services, among other risks.

"The Web is like the Wild West," says Roxana Geambasu, a computer scientist at Columbia Engineering and the Data Science Institute. "There's no oversight of how our data are being collected, exchanged and used."

With computer scientists, Augustin Chaintreau and Daniel Hsu, and graduate students Mathias Lecuyer, Riley Spahn and Yannis Spiliopoulos, Geambasu has designed a second-generation tool for bringing transparency to the Web. It's called Sunlight and builds on its predecessor, XRay, which linked ads shown to Gmail users with text in their emails, and recommendations on Amazon and YouTube with their shopping and viewing patterns. The researchers will present the new tool and a related study on Oct. 14 in Denver, at the Association for Computing Machinery's annual conference on security.

Sunlight works at a wider scale than XRay, and more accurately matches user-tailored ads and recommendations to tidbits of information supplied by users, the researchers say. Prior researchers have traced specific ads, product recommendations and prices to specific inputs like location, search terms and gender, one by one. One tool, AdFisher, received attention earlier this year after showing that fake Web users thought to be male job seekers were more likely than female job seekers to be shown ads for executive jobs when later visiting a news site.

| | email **subject** & text | ads **Title**, **url** & text | Results |
|---|---|---|---|
| **Race** | **Dominican**<br>dominican [...]  **(OR)**<br>**Hair**<br>hair cut hair cut | **Shampoo JC**<br>www.shampoojc.com<br>Professional Coloring, Highlights.<br>Make your appointment now! | p-value = 0.0004<br>1427 impressions<br>in 44 profiles<br>93% in context |
| **Religion & Spirituality** | **Mormon**<br>mormon mormon | **Family History Search**<br>genealogy.com/Family+History<br>1) Simply enter their name.<br>2) View their family history now! | p-value = 0.001<br>237 impressions<br>in 18 profiles<br>74% in context |
| | **Muslim**<br>muslim muslim | **Fine Models & Miniatures Shop**<br>www.1stdibs.com/Modern<br>Our Unique Modern Collection<br>Rare Items Added Every Week. | p-value = 0.007<br>59 impressions<br>in 11 profiles<br>100% in context |
| | **Jewish**<br>jewish jewish | **Free Ancestor Search**<br>archives.com<br>Looking for Your Family Ancestry?<br>Search for Free [...] | p-value = 0.0007<br>287 impressions<br>in 15 profiles<br>100% in context |
| | **Buddhist**<br>buddhist buddhist | **Berkshire Retreat**<br>www.eastover.com<br>Holistic Retreat Center Personal<br>Retreat | p-value = 0.008<br>190 impressions<br>in 17 profiles<br>43% in context |
| | **Guru**<br>guru spiritual guide  **(OR)**<br>**Astrology**<br>astrology psychic mystical | **What is Quantum Jumping?**<br>www.quantumjumping.com<br>Discover How Thousands of People<br>are Jumping to Change Their Life [...] | p-value = 0.001<br>244 impressions<br>In 34 profiles<br>76% in context |
| **Sexual Orient.** | **Gay**<br>gay homosexual<br>lesbian gay [...] | **Men Underwear/Workout**<br>www.gosoftwear.com<br>Underwear, Swimwear<br>Go Natural American [...] | p-value = 0.05<br>54 impressions<br>in 19 profiles<br>74% in context |
| **General Health** | **Affordable**<br>afforable care [...]  **(OR)**<br>**Nursing**<br>nursing home [...] | **Illinois Senior Living**<br>www.cottagesofnewlenox.com<br>Assisted Living for Seniors<br>in New Lenox [...] | p-value = 0.03<br>103 impressions<br>in 36 profiles<br>28% in context |
| | **Alzheimer**<br>Alzheimer Alzheimer | **1/3 of Seniors 65+ Fall**<br>jacuzzi-walk-in-tubs.com/Safety<br>Help Eliminate the Fear of Falling<br>in the Bathroom [...] | p-value = 0.01<br>21 impressions<br>in 8 profiles<br>100% in context |
| | **Depressed**<br>depression  **(OR)**<br>**Anxious**<br>anxious anxiety | **Is He A Cheater?**<br>spokeo.com/Cheating-Spouse-Search<br>Enter His Email Address. Find Pics &<br>Profiles From 70+ Social Networks. | p-value = 0.03<br>1179 impressions<br>in 52 profiles<br>20% in context |
| | **Cancer advice**<br>How did you cope with<br>cancer in your familly?<br>What an aweful disease! | **The Business of Wellness**<br>healthmediagroup.blogspot.com<br>What my doctor can learn from<br>my Shoe Shine Man [...] | p-value = 0.04<br>380 impressions<br>in 28 profiles<br>91% in context |
| **Prohibited** | **counterfeit, counterfeit**<br>counterfeit counterfeit | **A&A Global Industries**<br>aaglobal.com<br>Largest Supplier to Bulk Industry<br>Toys, Equipment, Candy, Supplies | p-value = 0.002<br>66 impressions<br>in 17 profiles<br>100% in context |
| | **drugs**<br>drugs cheap online order | **Eagle Creek Luggage**<br>www.eaglecreek.com/<br>Extremely Tough & Durable Gear.<br>Luggage, Organizers, Duffels & More | p-value = 0.03<br>214 impressions<br>in 19 profiles<br>99% in context |
| **Misc.** | **Deregulation**<br>deregulation [...] **(OR)**<br>**Financial Reform**<br>financial reform [...] | **Compliance Audit**<br>unifiedcompliance.com<br>Checklist All IT Compliance<br>You Need to Track In [...] | p-value = 0.0008<br>1582 impressions<br>in 36 accounts<br>61% in context |
| | **Unemployed**<br>lazy unemployed | **Easy Auto Financing**<br>www.midsouthautoloans.com<br>Need a quick car loan?<br>We work with credit issues | p-value = 0.006<br>161 impressions<br>in 24 profiles<br>8% in context |
| | **Payday**<br>payday loan | **Fast Cash Loan Online.**<br>www.checkintocash.com<br>Apply Now. Takes Only 5 Minutes.<br>It's as Easy as 1,2,3. | p-value = 0.007<br>198 impressions<br>in 10 profiles<br>6% in context |
| | **Veterans**<br>war veteran veterans | **Veterans Care Costa Rica**<br>www.veteranscarecostarica.com<br>Receive your proper medical care<br>Tricare, VA, Champ VA | p-value = 0.0006<br>490 impressions<br>in 15 profiles<br>84% in context |

Columbia researchers found evidence that some Gmail ads appeared to contradict Google's ban on using sensitive information in targeted ads. Credit: Columbia University

Sunlight, by contrast, is the first to analyze numerous inputs and outputs together to form hypotheses that are tested on a separate dataset carved out from the original. At the end, each hypothesis, and its linked input and output, is rated for statistical confidence. "We're trying to strike a balance between statistical confidence and scale so that we can start to see what's happening across the Web as a whole," said Hsu.

The researchers set up 119 Gmail accounts, and over a month last fall sent 300 messages with sensitive words in the subject line and body of the email. About 15 percent of the ads that followed appeared to be targeted; some seemed to contradict Google's policy to not target ads based "on race, religion, sexual orientation, health or sensitive financial categories," the researchers said. For example, words typed into the subject line of a message— "unemployed," "depressed," and "Jewish," were found to trigger ads for "easy auto financing," a service to find "cheating spouses," and a "free ancestor" search, respectively.

The researchers also set up fake browsing profiles and surfed the 40 most popular sites on the Web to see what ads popped up. They found that just 5 percent of the ads appeared to be targeted, but some seemed to violate Google's advertising ban on products and services facilitating drug use, they said. For example, a visit to "hightimes.com" triggered an ad for bongs at AquaLab Technologies, researchers said. Interestingly, the algorithms also seemed to pick up on the political leanings of popular news sites, pitching Israeli bonds to Fox News readers, and an anti-Tea

Party candidate to Huffington Post readers.

The researchers caution against inferring that Google and other companies are intentionally using sensitive information to target ads and recommendations. The flow of personal data on the Web has become so complex, they said, that companies themselves may not know how targeting is taking place.

In Nov. 10, 2014, Google abruptly shut down Gmail ads - the last day that Geambasu and her colleagues were able to collect data. The ads appear to have been replaced by so-called organic ads displayed in the promotions tab. Sunlight has the ability to detect targeting in those ads, too, said Geambasu, but the researchers haven't yet given that a try.

Sunlight's intended audience is regulators, consumer watchdogs and journalists. The tool lets them explore how personal information is being used and decide where closer investigation is needed, they said. "In many ways the Web has been a force for good, but there needs to be accountability if it's going to remain that way," said Chaintreau.

"Sunlight is distinctive in that it can examine multiple types of inputs simultaneously (e.g., gender, age, browsing activity) to develop hypotheses about which of these inputs impact certain outputs (e.g., ads on Gmail)," said Anupam Datta, a researcher at Carnegie Mellon who led the development of the AdFisher tool and was not involved in the current study. "This tool takes us closer to the critical goal of discovering personal data use effects at scale."

**More information:** www.cs.columbia.edu/~djhsu/papers/sunlight.pdf

Provided by Columbia University

Citation: New tool expands tracking of personal data on the Web (2015, October 9) retrieved 3 September 2024 from https://phys.org/news/2015-10-tool-tracking-personal-web.html