

The Evidentiary Value of Big Data Analysis

Marco Pollanen, Trent University, Canada
Bruce Cater, Trent University, Canada

The Asian Conference on Politics, Economics & Law 2016
Official Conference Proceedings

Abstract

Big data is transforming the way governments provide security to, and justice for, their citizens. It also has the potential to increase surveillance and government power. Geospecific information – from licence plate recognition and mobile phone data, biometric matches of DNA, facial recognition, financial transactions, and internet search history – is increasingly allowing government agencies to search and cross-reference. This heightened reliance on big data searches raises the question: what is the probative value of the information that results?

A distinguishing feature of the scientific method is that it begins with the development of an hypothesis that is then tested against data that either support or refute the hypothesis. That method is essentially followed in a conventional criminal investigation in which, after a suspect is first identified, evidence is gathered to then either build a case against, or rule out, that suspect.

The analysis of big data, however, can, at times, be more akin to first trawling for data to only later generate an hypothesis. In this paper, we investigate the conditions in which this may lead to problematic outcomes, such as more data leading to higher rates of false positives. We then sketch a big data analysis legal/policy framework that can circumvent these problems.

Keywords: database searches, forensic science, big data analysis, criminal databases

iafor

The International Academic Forum
www.iafor.org

Introduction

With the advent of smartphones, people may now leave behind them a near complete digital trail of their daily lives – from everything we search for and read online, to continuous location information, to a record of their interactions through social media and texting apps. Wearable technology is evolving in the direction of recording complete health and biometric information. Contactless transit cards, and arrays of video cameras coupled with improvements in facial recognition and license plate reading algorithms, allow our movements to be more accurately tracked. And cashless transactions allow every detail of purchases to be recorded.

At the same time, advances in computing and statistical analysis are increasingly allowing for the contents of these vast databases to be analyzed and for inferences to be drawn from those data.

In the United States, dragnets – where large numbers of people are indiscriminately questioned or detained – are a violation of civil liberties and considered unconstitutional. Despite this, recent years have seen growth in the area of digital dragnets that provide law enforcement with access to ever-larger DNA databases, financial and communications records, and the results of widespread electronic surveillance. These dragnets are often justified by threats to security or rationalized as resulting in only metadata. Regardless of how these searches are framed, this paper will show that they collect that may lead to erroneous conclusions.

An important, but often overlooked, aspect of the scientific method is that a hypothesis is formulated *before* data to test that hypothesis are drawn. Without establishing a hypothesis first, it is easy to fall into a trap of data dredging, in which inadvertent patterns are uncovered and misleading conclusions are drawn.

For that reason, modern forensic investigations ideally follow a path akin to the scientific method – a suspect is first identified (i.e., a hypothesis is first formed), then evidence to test that hypothesis is gathered.

In this paper, we will discuss the implications of deviating from that path. To introduce one type of problem that arises when hypotheses are not established first, the next section will outline a well-known result from probability theory known as the birthday paradox (Bloom, 1973). We will then illustrate something akin to the birthday paradox with real examples from DNA database matches. Finally, we will examine the situation with respect to Big Data searches and, and we will discuss potential solutions to the problems that arise.

The Birthday Paradox

Suppose that, in a group of N people, each of the 365 possible birthdays is equally probable; we will ignore leap years.

If $N = 2$, the probability that they share a birthday is $1/365 = 0.0027$ – that is, we could assign the first person any birthday, and the second person would have a $1/365$ chance of having the same birthday.

DNA Database Matches

In recent years there has been a lot of controversy about how to calculate the probability of a random DNA sample matching the profile of one found in a DNA database. While everyone has unique DNA, DNA databases typically only store a profile from this DNA consisting of measurements at a fixed set of locations (or loci) on the chromosome. Typically, 9 to 13 independent loci are selected for the database, with two unrelated samples matching at a particular loci with a probability of about 7.5%. Thus the odds two random unrelated profiles match at a fixed set of 9 loci is about 1 in 13 billion and at 13 loci is about 1 in 420 trillion.

One controversy that has arisen is from the statistical results from *The Arizona DNA Offender Database* (Kaye, 2009). At the time the database had 65,493 profiles, those profiles were analyzed and 122 pairs were found to match at 9 loci, 20 at 10 loci, and 1 pair at each 11 and 12 loci. Many people found these results astounding as the database was relatively small, and, as noted above, the probability of two random samples matching at 9 loci is about 1 in 13 billion, and at 12 loci about 1 in 32 trillion.

There are, however, several reasons why we would expect to see a large number of matches. The first is due to the birthday paradox, as described above. The second reason is that, in the case of 9 loci for example, the loci for which the matches occur could be different for different pairs of matches. From a set of 13 loci there are 715 different ways to choose 9 of them, so allowing partial matches increases the odds of a match by an additional factor of 715.

While these considerations do not fully explain the high number of matches, they do come close – for example, in the case of 9 loci, the expected number of matches would be 68, not 122 as were found. But, given the scale of the numbers being dealt with, that *is* fairly close, particularly given the crudeness of the genetic model in which it is assumed that all individuals are unrelated, and all loci are independent with equal probabilities of random matches. A more sophisticated analysis has been done by Mueller (2008).

One issue that arises immediately from DNA matches is that the science is relatively sophisticated and the odds of a random match can seem so overwhelmingly long that it seems possible to identify and convict a suspect by means of a DNA match only. But this is problematic if the match originated from a database search alone. A few cases are instructive.

In what was the first widely reported false match (Fowler, 2003) from a DNA database, a severely disabled man in the United Kingdom who was arrested for a burglary that occurred some 200 miles away and that involved the burglar climbing through a window. In that case, the only evidence was a match from a database search with a probability of 1 in 37 million, which corresponds to 6 loci. The near impossibility of the man committing the crime did not clear him.

With the population of the United Kingdom being 64 million, on average we would expect any 6-loci DNA profile to be shared by two people. But, by conducting a DNA database search, we essentially trawled through millions of hypotheses to fit the evidence, violating the first tenant of the scientific method – that we must first have a

hypothesis. This illustrates what is known as the prosecutor's fallacy (Thompson & Shumann, 1987), in which investigations and prosecutions revolve around a probability of match. The correct interpretation is that if the suspect is innocent, there is a 1 in 37 million chance that there is a match. However, with the prosecutor's fallacy, the clauses are reversed and the logically incorrect interpretation is adopted – if the DNA matches, there is a 1 in 37 million chance that the suspect is innocent.

It is not just investigators and prosecutors who incorrectly weigh DNA evidence. A 30-year old cold-case (Murphy, 2015) facilitated the analysis of partial matches in the Arizona DNA database. The defendant in that case was identified and convicted largely due to the partial match of the badly degraded DNA sample to a profile found in a California database. The judge allowed only the prosecution's statistic that the chance that an individual picked at random would match the crime-scene DNA is 1 in 1.1 million. Jurors were not informed that the match was a result of a database trawl, whereby 9-loci partial matches are not uncommon, nor were they informed that about 40 people in California would be expected to have a profile that matches the crime-scene sample. The fact that a partial match was used is not that uncommon, as crime scene evidence can be degraded and mixtures of DNA samples can result. Furthermore, different databases often use different loci for profiles, and searches can be done using the profiles of close relatives.

A further problem with assigning astronomical probabilities to a single piece of evidence, such as a DNA database match, is that those probabilities would be dwarfed by real-life considerations, such as laboratory errors and contamination. For example, a man in Australia was convicted of raping a woman found unconscious at a nightclub based solely on a random match in the Australian DNA database (Roberts & Hunter, 2012), despite other evidence suggesting that the individual could not be a suspect. Only through post-conviction serendipity was it discovered that the original rape-kit was likely contaminated at the laboratory, leaving no clear evidence that a crime even took place.

Even when evidence is found at the crime-scene and it is correctly attributed to an individual, the relevance of the sample to the crime must be established. Typically, DNA establishes, at most, the presence of or contact with an individual, not that they committed a crime. In another case, a man in the United Kingdom (Barnes, 2012) was jailed for eight months when a partial match was found between his DNA profile in a database and a crime-scene sample from a murder scene. It has been suggested that because the suspect was a taxi driver, he likely came into contact with the victim individual and some of his shed skin cells clung to that person.

While there are potential pitfalls in interpreting DNA evidence, especially when it comes from random matches found by trawling through databases, it is important to note that DNA evidence is still some of the most reliable types of evidence there is, and that it has likely lead to the exoneration of more people far more often than it has resulted in false convictions. By comparison, while identification by eyewitnesses carries a lot of weight in courts, studies have shown how utterly unreliable eyewitness testimony can be (National Research Council Report, 2014). We introduced the issues with the *Arizona DNA Offender Database* to demonstrate how the birthday problem arises in criminal investigations.

In the next section we will discuss how these problems might be amplified as the number and type of databases used in forensic investigations increases.

Big Data Searches

In recent years, aided by technological advances and often rationalized as necessary to fight terrorism, mass surveillance has been increasing. For example, in the United States, metadata for hundreds of billions of telephone calls has been collected (Cauley, 2006); the exterior of all letter mail is photographed (Miga, 2013); databases containing information on financial transactions, e-mails, and internet surfing habits are maintained; and social media are monitored (Kawamoto, 2006). The FBI has a face-recognition system with a database of over 400 million photos (Kelley, 2016). Combined with the ever increasing array of CCTV cameras, it might be possible to recognize individuals in any public location. For example, the United Kingdom has up to 6 million CCTV cameras (Barrett, 2013), about one for every 11 individuals. Furthermore, location information could also be obtained from license plate recognition or from databases of transit card usage.

In addition to government databases, private companies such as Google and Facebook have access to vast amounts of information about individuals, except where one makes exerts considerable effort to maintain their privacy. This is especially true due to the near ubiquitous use of smartphones. They potentially have access to the contents of every digital communication one partakes in, and to one's location history; they can map one's photographs, social connections, and browsing and search histories; and they can potentially track health and biometric information through a phone's sensors. This information is also available to governments seeking to increase surveillance.

The average individual leaves a vast digital trail throughout the day, from which it may be possible to surmise when he/she woke up and how long they slept, their location throughout the day, including where they work, where they shop, and what they bought, read, or wrote. By combining the available information with information from biometric sensors in smartphones or wearable devices, it may be possible to develop algorithms that give an idea of what an individual thought and felt throughout the day or to predict behavior.

While all of this information can be a boon for law enforcement in their quest to solve crimes, as the number and size of databases grow it also has the potential to lead to an increase in the number of falsely accused individuals. To see this, consider a DNA database in which an individual profile will typically contain information regarding 13 loci. This would be equivalent to an individual having a record in 13 different databases, each containing the information of a single loci. Thus searching through multiple databases of digital information would also be subject to the Birthday Paradox as we have seen with DNA. Moreover, when an individual matches information in only some databases but not others, this further magnifies the problem of false identifications from partial DNA matches have been shown to have with the Arizona DNA database.

While there are many similarities with searching through digital information databases and DNA databases, there are causes for greater concern. DNA analysis

occurs in a laboratory setting, and while the measurements have errors associated with them, they can be estimated. Laboratory errors do occur, of course, but it is still a scientific setting where one would believe every attempt would be made to estimate and minimize these errors. On the other hand, analysis of databases of other digital records might involve information that was not originally intended for forensic examination, such as facial recognition on grainy photos or the inaccuracies of finding the location of a mobile phone user. These errors may be poorly understood and might contribute significantly to birthday paradox collisions. Furthermore, DNA analysis involves trying to match a set number of loci, while a trawl through digital data may involve an unknown number of databases and is problematic because the probability of a match would be incalculable. We have seen that partial matches of only some loci significantly increase the chance of misidentification with DNA databases, but in that case we know which loci cannot be matched and probabilities could be adjusted accordingly. It would be even more problematic if the databases investigated were not known or revealed. For example, suppose while in an investigation, all Google searches for a particular explosive were flagged. If a suspect had searched for that particular compound, that would certainly be used to build the case against him or her. On the other hand, if the suspect was *not* one of the individuals who had made that particular search, that fact might not be factored into the calculation of their probable guilt and it would almost certainly be inadmissible in court.

The concern with searching through a large number of databases for suspects that could fit the evidence of a crime is, of course, that people may be falsely accused. But with such an overwhelming amount of circumstantial evidence pointing to them, it could be difficult to exonerate them. For example, perhaps a murder has been committed, and by pure chance alone an individual is found whose license plate was caught driving nearby at a similar time, traces of whose DNA are found on the murder victim (perhaps because they ate at the same restaurant), and perhaps the day before they bought the same brand of duct tape used in the crime. As technology and pattern recognition algorithms get better, it is likely that even more casual links in vast arrays of data will be found.

Discussion/Conclusion

Databases are important tools for fighting crime and protecting national security. Indeed, as criminals become more sophisticated in their use of technology, there is arguably a need for law enforcement to do the same. A significant problem arises, however, because trawling through a database without a suspect in mind – essentially in violation of the scientific method – may result in erroneous conclusions.

A primary goal, then, should be to put these searches on a more scientific footing. Perhaps the most obvious way to achieve this would be to use the database for identifying a suspect (i.e., formulating a hypothesis), and not for building a case against him or her. Only further evidence gathered from other sources would be used for the purposes of prosecuting. A second possible approach would involve (perhaps randomly) separating a list of all available databases into two parts. From one part, a suspect could be identified, while, from the second part, searches could be conducted to build a case against that suspect. Of course, many would object to either approach, for they would be seen as leaving evidence unused.

Whether or not a hypothesis is found first, it is important to understand the statistical characteristics of many of the key databases used in order to understand their scope and potential inaccuracies. This would be important for assigning a probability of a match for use in the legal system.

Of course, even in the case of DNA, law enforcement agencies have fought access by researchers to study random match probabilities (Kaye, 2009). Yet, potential violations of privacy could be circumvented by removing any to individuals. The data could be scrambled or encrypted in some way without compromising researchers' abilities to analyze the key characteristics of the data. In order to make investigations more scientific, it is important to carefully document all database searches included in the hunt for a suspect, even the ones that lead to negative results. For these purposes, the development of a standardized set of databases and search criteria would be appropriate.

References

Barrett, D. (2013, July 10). One surveillance camera for every 11 people in Britain, says CCTV survey. *The Telegraph*. Retrieved September 7, 2016, from <http://www.telegraph.co.uk/technology/10172298/One-surveillance-camera-for-every-11-people-in-Britain-says-CCTV-survey.html>

Barnes, H. (2012, August 31). DNA test jailed innocent man for murder. *BBC News*. Retrieved September 6, 2016, from <http://www.bbc.com/news/science-environment-19412819>

Bloom, D. (1973). "A Birthday Problem". *American Mathematical Monthly*. 80: 1141–1142.

Cauley, L. (2006, May 11). Advertisement NSA has massive database of Americans' phone calls. *USA Today*. Retrieved September 7, 2016, from http://usatoday30.usatoday.com/news/washington/2006-05-10-nsa_x.htm

Fowler, R. (2003, April 27). DNA, the second revolution. *The Guardian*. Retrieved September 5, 2016, from <https://www.theguardian.com/uk/2003/apr/27/ukcrime7>

Kaye, D. H. (2009). Trawling DNA Databases for Partial Matches: What Is the FBI Afraid Of?. *Cornell Journal of Law and Public Policy*, 19(1).

Kawamoto, D. (2006, June 9). Is the NSA reading your MySpace profile? *CNET*. Retrieved September 7, 2016, from http://archive.is/20120720043006/http://news.com.com/2061-10789_3-6082047.html#selection-925.5-929.1

Kelly, H. (2016, June 16). FBI's face-recognition system searches 411 million photos. *CNN Money*. Retrieved September 7, 2016, from http://money.cnn.com/2016/06/16/technology/fbi-facial-recognition/index.html?iid=ob_homepage_tech_pool

Miga, A. A. (2013, August 2). AP Interview: USPS takes photos of all mail. *AP Online*. Retrieved September 7, 2016, from http://www.highbeam.com/doc/1A1-10eae68abcc8439fa78610fe561ab6fc.html?refid=easy_hf

Mueller, L. D. (2008). Can simple population genetic models reconcile partial match frequencies observed in large forensic databases?. *Journal of genetics*, 87(2), 101-108.

Murphy, E. E. (2015, October 08). The Dark Side of DNA Databases. *The Atlantic*. Retrieved September 6, 2016, from <http://www.theatlantic.com/science/archive/2015/10/the-dark-side-of-dna-databases/408709/>

National Research Council of the National Academies. (2014). *Identifying the culprit: Assessing eyewitness identification*. Washington, D.C.: National Academies Press.

Roberts, P., & Hunter, J. (2012, May 18). *Criminal evidence and human rights: Reimagining common law procedural traditions* (P. Roberts & J. Hunter). Bloomsbury Publishing.

Thompson, E.L.; Shumann, E. L. (1987). "Interpretation of Statistical Evidence in Criminal Trials: The Prosecutor's Fallacy and the Defense Attorney's Fallacy". *Law and Human Behavior*. 2 (3): 167. doi:10.1007/BF01044641.

Contact email: marcopollanen@trentu.ca