

WM0824TU Deliverable 1 - Security metrics

Adriaan de Vos, Bas van 't Spijker, Djoshua Moonen, Kes Greuter

(a.d.devos@student.tudelft.nl, b.c.vantspijker@student.utwente.nl, D.D.M.Moonen@student.tudelft.nl,
k.o.greuter@student.utwente.nl)

September 29, 2020

1. What security issue does the data speak to?

The goal of this deliverable is to analyze possible metrics and their respective performance that can be created from the Privacy Rights Data Breaches data set [4]. This data set contains information about 9015 US data breaches that have been reported to state Attorneys General and the U.S. Department of Health and Human Services. Data breaches can be defined as a security violation in which sensitive, protected or confidential data is copied, transmitted, viewed, stolen or used by an unauthorized individual [5]. The given data set contains information about the location, date, number of records, type of breach, and the type of organisation.

According to Privacy Rights, all states in America are legally obligated to report a data breach. There are, however, significant differences in the definitions regarding data breaches, such as the definitions of personal data, what constitutes a breach, the requirements for a notice and exemptions [3]. The differences in the regulations per state cause skews in the data set, as the states report different information and have different standards to when a data breach should be reported in the first place. The analysis of the metrics within the database should therefore take these differences in account by not focusing purely on data about in which areas some states lack, or other similar possible biases in the database.

For the US government it is important to have an overview of the performance of a certain sector. To do achieve this, the frequency of data breaches and the financial impact that a data breach has on a sector are two factors that can be determined by making use of this database. A data breach of personal data can not only result in reputation damage, but competitors could possibly also take great advantage of the data leaked. Also, customers could become a target for identity theft or fraud. Organizations thus want their sensitive data to be secure against cyber attacks [15]. Therefore, it is useful for the US government to see what still causes such data breaches. The actor in our security metrics then, is the US government, that has interest in the performance of different sectors regarding data breaches. The security issue that will be regarded in this assignment is whether there are big differences in impact and frequency of data breaches in a sector and what the possible causes of this might be.

2. What would be the ideal metrics for security decision makers?

The ideal metric for the US government regarding data breaches would be the total financial impact and frequency of data breaches on the different sectors. Information on the financial impact would indicate whether it is appropriate for these sectors to invest more in security taking into account the losses from cyber security breaches. According to the Gorden-and-Loeb model for example, there should not be more than 37 percent of the expected losses from cyber security breaches spent on improvement of security [10].

The financial impact of a data breach of a sector includes direct and indirect costs and should also be put in proportion with the number of attempts of data breaches on the sectors. One of these direct costs is the governmental fines. One thing to note with this cost is that regulations and perhaps also the cost of the fines differ per state, thus these costs could differ when one analyzes the costs per sector per state. Other direct costs include the costs of possible lawsuits due to customer damage and the costs of losing customer subscriptions. Indirect costs that might occur due to a data breach are the costs of possible advantages for

obtained by competitors and the costs of inconvenience for employees due to e.g. reconsidering requirements and reorganization.

The ideal metric would include the impact of each data breach per sector by adding the above direct and indirect costs and multiplying these costs times the frequency of such a data breach within a sector.

$$I_s = F * I_F \quad (1)$$

The financial impact I_F is further divided into internal and external costs. The internal costs come from actors within the sector such as repairing sector image. The external costs come from potential fines and lawsuits. Formally this is written down as:

$$I_F = C_e + C_i \quad (2)$$

Equation 1 includes information that is not available in the database, and perhaps not even available at all. Besides that, in case data on the value of loss is available, this value might not be completely accurate as companies often inflate costs, or base the costs on weak evidence or simplified assumptions [1]. However, with the available data an estimate of the impact of a data breach can still be made by using the number of records loss per data breach and the frequency of such data breaches per sector. The records loss could give an indication of the number of customers that suffered damage, which in turn increases the sector impact. Yet, as there are no monetary values for the personal data of end-users, there still needs to be an estimate made of the costs per customer that suffered damage [16]. To make a precise estimation there should be more empirical data available to allow for calculation of a more evidence-based metric. The lack of publicly available cyber security data sets is currently a limiting factor within the research field [17].

3. What are the metrics that exist in practice?

Metrics can be based on measuring either controls, vulnerabilities, incidents, or (prevented) losses [2]. Controls are measures with the purpose of mitigating risk. Vulnerabilities are the metrics that evaluate the performance of these controls given a certain threat scenario. The exploitation of such a vulnerability leads to an incident. In a case of a vulnerability, adversaries have attempted to attack the controls in place. There are also metrics that map incident on (prevented) losses. The database used in this assignment gives an overview of data breach incidents.

Figure 1 shows the focuses four commonly known methods to organize metrics. From this scheme, it can be concluded that all of these methods have a focus towards controls and vulnerabilities, rather than incidents and none of these metrics have a focus towards losses. Methods benefit from focusing towards controls, as measuring controls is easier than measuring vulnerabilities[2], incidents and losses. Also, an organization can use methods about controls as a way to improve the sale figures of security solutions. Finally, the focus on controls leaves the risk of failure at the buyer rather than the organisation.

The Privacy Breach database [4] includes information on incidents by the reporting of data breaches. The number of records however, is a value that could be used in the calculation of losses.

The data set solely gives information on incidents, as it includes only the type of breach, the organisation that held the records, and the number of records breached. Although this data does not include incident metrics by itself, the data can provide threat level indicators, such as the frequency of an successful attack per type of organisation. Such an indicator alone is not precise enough for a metric [13], yet it could help update quantitative risk assessment schemes and inform a type of organisation on the most common type of breach for such an organization, and thus influence the defense tactics of that organisation.

4. Which metrics can be designed from the data set?

	Controls	Vulnerabilities	Incidents	(Prevented) losses
Cloud Security Alliance Cloud Control Matrix				
Security Service Level Agreement (SLA)				
Software Security Maturity Model				
Cyber Security Assessment Netherlands				

Covers fully	
Covers largely	
Covers partly	
Covers slightly	
Does not cover	

Figure 1: Comparison of metrics

The database contains several columns of data, of which not every column is equally useful. The 'date made public' is not very helpful since there is no 'date of discovery' to compare it to. There are 7624 rows with a unique 'company name', but the most reported breaches a single company has had, is 13. The database also has city, state, and longitude/latitude information per breach. Localization information is most likely not the most important, given that the online landscape (locations of servers, amount of data stored) is not necessarily congruent to the physical landscape, and state legislation on privacy breaches differs greatly[3]. Still, it could be interesting to take a look at how often breaches were reported per state. The most interesting columns however, are 'type of breach' (card fraud, hack, insider job, physical documents compromised, stolen portable device, stationary computer/server loss, unintended disclosure, unknown), 'Total records', 'Type of organisation' (BSF: financial services, BSR: retail, BSO: other businesses, educational institutions, government and military, medical, nonprofits, unknown) and 'Year of breach.' Based on these, seven metrics are designed. It has to be noted that 2187 out of the 9015 entries list '0' at 'total records lost'. The reasons vary: sometimes the exact number of records is unknown, and sometimes the number is very low (e.g. a single employee lost their ID card). Since checking every '0' row for legitimacy is out of the scope of this assignment, it is assumed that all '0' rows are valid.

The first metric is a graph of the total number of breaches per year. This will give a general overview of how the threat environment is developing. The formula for the cardinality of each subset S per year y , consisting out of rows r , out of the complete dataset D , can be denoted as:

$$|S_y|, \text{ where } S_y = \{r \in D | r_{year} = y\} \quad (3)$$

The next metric is a graph of the number of breaches per type over the years. This will give an indication of the type of breaches that might be getting more prevalent. The formula for the cardinality of each subset S per year y and type t , consisting of rows r , out of the complete dataset D , is given as:

$$|S_{y,t}|, \text{ where } S_{y,t} = \{r \in D | r_{year} = y \wedge r_{type} = t\} \quad (4)$$

Knowing what type of breach is the most likely to occur is useful, but it is also desirable to know which sectors have been hit the hardest.

Therefore, the next metric is similar to the second, but instead of breaches per type over the years, we calculate the breaches per sector over the years. A similar formula as the previous one is given for the calculation:

$$|S_{y,s}|, \text{ where } S_{y,s} = \{r \in D | r_{year} = y \wedge r_{sector} = s\} \quad (5)$$

Large dataleaks seem to be occurring more, so it would be interesting to see if the dataset reflects that. That is why the next metric shows the total number of records leaked per sector per year. This will show the impact of breaches in general over the years, and which sectors leak the most data. The sum of records per sector per year is given by:

$$\sum_{r \in S_{y,s}} r_{records}, \text{ where } S_{y,s} = \{r \in D | r_{year} = y \wedge r_{sector} = s\} \quad (6)$$

We are also interested in seeing if, on average, there have been more records leaked per breach. This might show whether breaches have become more dangerous. Therefore, the fifth metric designed is the average number of records leaked per breach per sector per year. The formula is attained by dividing metric 6 by metric 5:

$$\frac{\sum_{r \in S_{y,s}} r_{breach}}{|S_{y,s}|}, \text{ where } S_{y,s} = \{r \in D | r_{year} = y \wedge r_{sector} = s\} \quad (7)$$

To investigate how localization might influence privacy breaches, we can look at the number of breaches or the number of records lost, either per state or per city. Since the distribution of records lost is heavily skewed (see section 5), number of breaches seems more useful. A city scale would be a bit too small for this research, so state scale is taken. The formula for the cardinality of S per state in database D is then given by:

$$|S_t|, \text{ where } S_t = \{r \in D | r_{state} = t\} \quad (8)$$

Lastly, one would expect there to be more breaches in a state where more people live. Therefore, the last metric is the amount of people per breach per state. This metric is calculated by dividing population P of state t by the cardinality of S :

$$\frac{P_t}{|S_t|}, \text{ where } S_t = \{r \in D | r_{state} = t\} \quad (9)$$

5. Evaluation of the metrics you have designed.

From figure 2 to 4 it becomes somewhat clear that the data is incomplete. After 2016, the number of breaches starts to become less. This however, is most likely due to the fact that it can take years for a breach to get officially reported. 2019 and 2020 are already not taken into account, since the number of (reported) breaches in those years are close to zero in the dataset. Otherwise, from figure 2 we can see that the number of reported breaches has stayed roughly the same since its last major increase in 2010, to around 750 breaches per year. Metric 3 shows that the number of hacks have been on the rise since 2009, as the number of insider-related compromises, together with the number of times that physical documents were stolen or portable devices lost, have been decreasing. The unintended disclosure number has somewhat increased since 2009 as well. This might, among other things, have to do with the popularity of doxxing [7]. Figure 4 shows that in terms of sheer number of breaches, the medical sector has been hit the hardest. Second to that comes the other businesses (BSO in the graph).

Figure 5 depicts the total number of records leaked in data breaches for each of the different sectors. It shows that the BSO sector took an enormous hit in 2016 and has been recovering years after. This is mostly due to

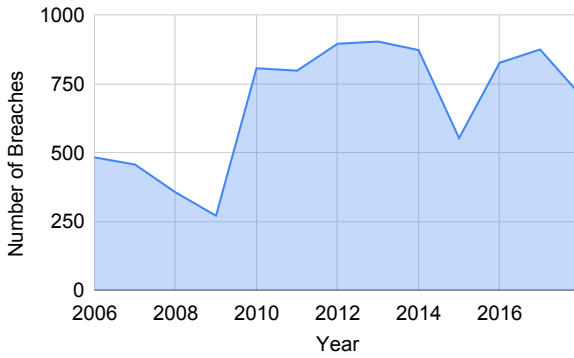


Figure 2: Metric #1: Number of Breaches per Year

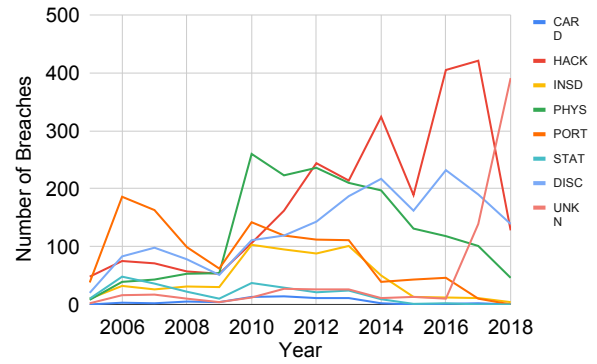


Figure 3: Metric #2: Number of Breaches per Type over the years

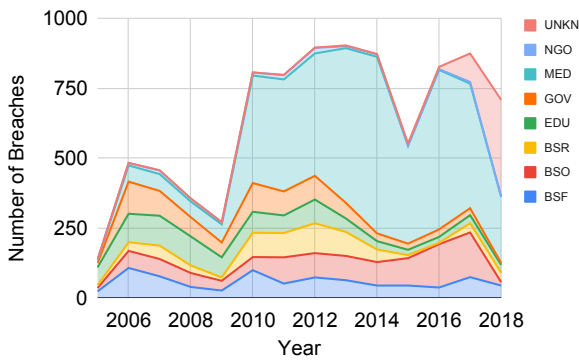


Figure 4: Metric #3: Number of Breaches per Sector over the years

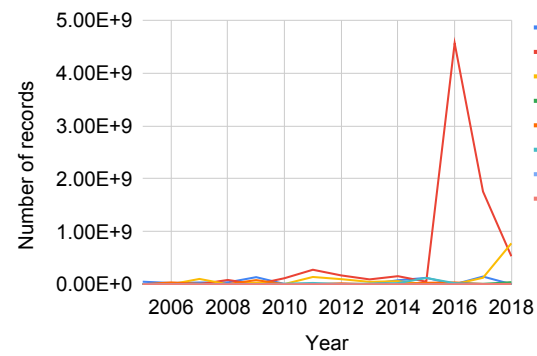


Figure 5: Metric #4: Total Records leaked per Sector over the years

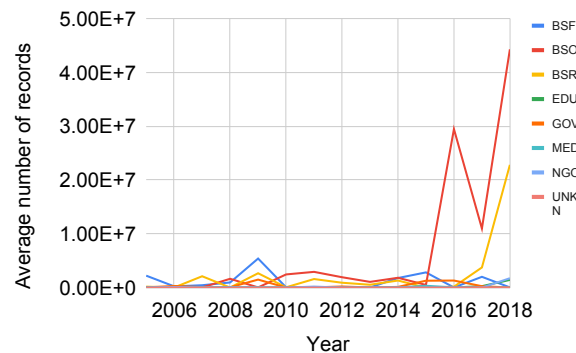


Figure 6: Metric #5: Average Records Leaked per Breach Per Sector over the years

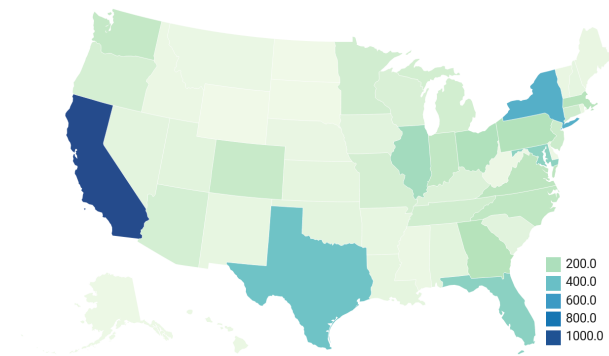


Figure 7: Metric #6: Total Breaches per State

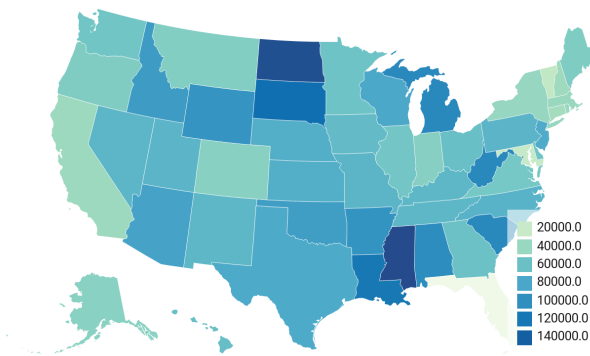


Figure 8: Metric #7: Population per Breach

a leak of 3 billion records at Yahoo, but although the breach was discovered in 2016, the actual breach year was 2013[9]. There were also some large breaches of a hundred million at several social media sites, such as LinkedIn, that year. Similar but less obvious patterns can be found in other sectors where the total number increases in one year, and decreases over the few years after that.

It is clear that this metric is impacted heavily by data breaches where many records are leaked. At several points, through a couple major data leaks, the total number of records leaked of the entire sector spikes.

Figure 6 shows average number of leaked records per sector per breach. Because the number of leaked records differs vastly per sector, and because of the large data spikes, it is hard to pull any meaningful conclusions from this metric. This is reinforced by the fact that, when not taking into account the large spikes from 2016 to 2018, the numbers vary wildly per year for every single sector.

Looking at figure 7, it might be assumed that the state of California is hit disproportionately hard compared to the rest of the U.S., but figure 8 disputes this view by correcting for population; in fact, California has relatively little breaches, and the states with less people have relatively more breaches. However, the map in general becomes much more homogenized. From these two maps we can conclude that although the number of breaches increases with population, it does so at a declining rate. It has to be kept in mind though, that state legislation and procedures might significantly influence the numbers per state[3].

6. What are the relevant differences in security performance observed?

The security performance implies how resistant a sector is against the security issue as described in Q1. The comparison of the number of hacks succeeded in different sectors could provide information on the quality of the security measures in place in these sectors. However, to accurately make an estimate of the resistance against breaches of a sector, one must have additional information such as the number of attempts performed at this sector and the size of this sector. When comparing the number of breaches per sector over time, as can be seen in figure 4, it can be concluded that the medical (MED) sector experienced at least three times more breaches than other sectors. Additional information on the MED sector is needed to conclude whether the security of the medical sector is of less quality than other sectors. Research shows that there has been a trend [11] in attacking the MED sector, as they store personal medical information about clients. While the MED sector is less frequently attacked than the financial sector [12], it appears that the most privacy breaches occur in the MED sector. Due to a lack of adequate IT spending and a lack of awareness of cyber crime, vulnerabilities within the MED sector could be exploited[6].

7. What risk strategies can the problem owner follow?

The US government can use the insight provided by our designed security metrics to govern individual states, improve current regulations within industry sectors, or make nation-wide efforts to reduce the overall risk. We have shown that there exists significant differences between state regulations [3], and we have also shown that there are significant differences between industry sectors in figure 4.

It would be best for the problem owner to make the first and foremost priority to create better insight in the impact of data breach incidents. The current differences that can be seen between industry sectors can partially be explained by different definitions of personal data and the different requirements for reporting data breaches. Ideally, the data available should show a very precise and realistic overview of all data breaches within the US. This could be achieved by extending and normalizing the current state regulations. A more precise overview would allow government to make better informed decisions about the risk strategies to follow.

Secondly, there should be a focus on reducing the risks of certain industry sectors that are already shown, by the current dataset, to be hit hard. It can be seen from figure 6 that the BSO industry sector (online companies) should probably be the first one that the government needs to look into because of the high number of records leaked per data breach. Alternatively, figure 4 shows that there are also many data breaches within the medical industry sector. We would expect the risk acceptance to be lower for the medical sector because of the very sensitive private medical data that is leaked. Therefore the US government should also focus on this sector when trying to reduce the risk.

Reducing the risk caused by data breaches should be very useful for the US government. After all, they are expected to protect both the citizens and the economy of the US against bad actors. There are already some examples of risk strategies that reduce leaking private data currently used by the US government [14] and the Indian government [8] to protect the citizens against foreign influence.

8. References

- [1] Hadi Asghari, Michel van Eeten, and Johannes M Bauer. "Economics of cybersecurity". In: *Handbook on the Economics of the Internet*. Edward Elgar Publishing, 2016.
- [2] TU Delft Carlos Gañán. *Economics of cybersecurity - Metrics in practice*. Available at <https://www.youtube.com/watch?v=GuyKbRoLmcw> (January 25th, 2015).
- [3] Privacy Rights Clearinghouse. *Data Breach Notification in the United States and Territories*. Available at <https://privacyrights.org/resources/data-breach-notification-united-states-and-territories> (December 10th, 2018).
- [4] Privacy Rights Clearinghouse. *The Chronology of Data Breaches!?* Available at <https://privacyrights.org/data-breaches> (September 19th, 2020).
- [5] Privacy Rights Clearinghouse. *What's a data breach?* Available at <https://privacyrights.org/resources/whats-data-breach> (July 11th, 2019).
- [6] The Hague Security Delta. *Healthcare Sector Report: Cyber Security for the Healthcare Sector*. Available at https://www.thehaguesecuritydelta.com/media/com_hsd/report/291/document/5ad7266dc1cba.pdf (2018).
- [7] Stine Eckert and Jade Metzger-Riftkin. "Doxxing". In: *The International Encyclopedia of Gender, Media, and Communication*. American Cancer Society, 2020. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/9781119429128.iegmc009>. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781119429128.iegmc009>.
- [8] Ministry of Electronics IT. *Government Blocks 118 Mobile Apps Which are Prejudicial to Sovereignty and Integrity of India, Defence of India, Security of State and Public Order*. Available at <https://pib.gov.in/PressReleasePage.aspx?PRID=1650669> (September 2nd, 2020).

- [9] Jim Finkle and Anya George Tharakan. *Yahoo says one billion accounts exposed in newly discovered security breach*. 2016. URL: <https://www.reuters.com/article/us-yahoo-cyber-idUSKBN1432WZ> (visited on 09/28/2020).
- [10] Lawrence A Gordon and Martin P Loeb. “The economics of information security investment”. In: *ACM Transactions on Information and System Security (TISSEC)* 5.4 (2002), pp. 438–457.
- [11] Zhang G Hu Y Ge L. “Research on differential privacy for medical health big data processing”. In: *20th International Conference on Parallel and Distributed Computing, Applications and Technologies* 20 (2019), pp. 140–145.
- [12] KPMG. *Health Care and Cyber Security: Increasing Threats Require Increased Capabilities*. Available at <https://assets.kpmg/content/dam/kpmg/pdf/2015/09/cyber-health-care-survey-kpmg-2015.pdf> (2015).
- [13] Böhme R. “Security Metrics and Security Investment Models”. In: *Advances in Information and Computer Security* 6434 (2010).
- [14] Wilbur Ross. *Commerce Department Prohibits WeChat and TikTok Transactions to Protect the National Security of the United States*. Available at <https://www.commerce.gov/news/press-releases/2020/09/commerce-department-prohibits-wechat-and-tiktok-transactions-protect> (September 18th, 2020).
- [15] Daniel J Solove and Danielle Keats Citron. “Risk and anxiety: A theory of data-breach harms”. In: *Tex. L. Rev.* 96 (2017), p. 737.
- [16] Hai Tao et al. “Economic perspective analysis of protecting big data security and privacy”. In: *Future Generation Computer Systems* 98 (2019), pp. 660–671.
- [17] Muwei Zheng et al. “Cybersecurity research datasets: taxonomy and empirical analysis”. In: *11th {USENIX} Workshop on Cyber Security Experimentation and Test ({CSET} 18)*. 2018.