**COMPARISON AND EVALUATION OF USER PRIVACY RISK IN RECOMMENDATION SYSTEMS**

A Research Paper submitted to the Department of Engineering and Society
In Partial Fulfillment of the Requirements for the Degree
Bachelor of Science in Computer Science

By

Shivaen Ramshetty

March 25, 2021

On my honor as a University student, I have neither given nor received unauthorized aid on this assignment as defined by the Honor Guidelines for Thesis-Related Assignments.

ADVISOR
Catherine D. Baritaud, Department of Engineering and Society

**ILLUMINATING THE DANGERS OF DATA COLLECTION AND STORAGE IN RECOMMENDATION SYSTEMS**

The recent impressive growth of the recommender systems has profoundly shaped the online marketplace. Currently, recent advancements in the recommendation system arena have brought attention to the industry, starting with and most notably the Netflix Prize. The Netflix Prize was a one-million-dollar award going to the group that created the best recommendation algorithm for Netflix's content recommendation engine (Bennet et al., 2007, p. 1). Consequently, such a large competition and all the positive interest that came, created massive funding for research and system upgrades. Companies like Google and Amazon made improvements to their recommendation systems, with Google upgrading their engine in 2016 to one powered by Google Brain (artificial intelligence) (Faggella, 2017). From Amazon's product recommendations to Netflix's "Top Picks for You," predicting user preferences or offering the next suggestion is growing more and more paramount, since the potential profit gained from each individual user grows when their recommendations more accurately assess their needs/wants. Wider application of this system has translated into the fact that "30% of Amazon's page views result from recommendations" and "80% of the content watched by Netflix subscribers comes through personalized recommendations" (Adomavicius et al., 2018, p. 2). However, the information these systems collect presents a growing amount of risk for users due to the quantity and nature of what is collected.

While the technical work addresses the problem of popularity bias in recommendation models themselves, where recommendations are reflective of most popular items rather than the best item, the STS work focuses on the risk users face through the collection of their sensitive data and the vulnerabilities that arise. With the previously mentioned growth of recommender

systems comes the fear that their information consumption may lead to greedy/dangerous activities from both the companies operating the system and external agents. These activities include leaks, collection of more data, collection of more private data, hacks, and the models that the technical work remedies. However, it is not clear whether recommendation systems pose any more significant a threat to users than other online systems, such as notification systems. To come to a resolution on this matter, this paper will employ a technology and social relationships research approach alongside Mesthene's technical complexity model to define risk and the qualities of a system that users find appealing in terms of privacy (Mesthene, 1970). Finally, by comparing the differences between recommendation systems and other online systems the research will explore what steps need to be taken and the ways in which solutions can be implemented.

**RISK ANALYSIS ON RECOMMENDATIONS TO FIND AREAS OF VULNERABILITY**

Recommendation systems face an "ethical dilemma", as Martin and Schinzinger (2010) coin, which is the situation where the benefit provided by applications making use of recommendation systems comes at the cost of user privacy (p. 27). The benefits of such a system can be listed as: its wide accessibility, customer retention, customer engagement, and customer satisfaction. Each one of these advantages is mutually beneficial for the company utilizing the system and also the consumer base. In general, the more aligned recommendations are to the user's preferences the more money the system can earn for the company. Consequently, managing and balancing this trade-off between consumer utility and the cost of their privacy is the main challenge for all actors within this space.

Referring to the benefits that recommendations provide, Jeckmans et al. (2013) state that, "recommender systems can meet the demands of large online applications that operate on a

global scale" (p. 2). As a result, corporations are extremely cognizant of the need to profile/group

similar users and understand how to meet the needs of such a variety of people. In fact, users

themselves know that "sites commonly track their browsing patterns, purchase histories, and

other sources of data to present individually personalized suggestions" (Harley, 2018, para. 5).

Therefore, data's extreme versatility allows entities to assess and serve all their consumers, but

its prevalence in current economic and social domains renders it a point of contention due to the

manners it can be mishandled.

**THE NEED TO BETTER UNDERSTAND DATA'S IMPACT FOR RISK ANALYSIS**

Data has been thrust into the forefront of the online ecosystem in recent years due to the

monetary value it holds for large companies. Provided by users willingly or not, corporations all

over the world use it to improve user experience within an application or to buff profits. The real

problem is whether or not the system poses an inherent risk to users through its data collection

practices. To define risk in this scenario, it is first necessary to understand what privacy and

confidentiality are in online systems. For the former, Jeckmans et al. (2013) introduces it as ""an

individual's claim to control the terms [of their] personal information"," while defining the latter

as "secrecy of individual pieces of information" (p. 6-7).

Now, with these terms in mind the discussion of risks and how they affect the system is

possible. Risks come from the vulnerability of hacks, the type of data stored, the possibility for

system exploitation, among many other sources. As Figure 1 on page 4 illustrates, these

vulnerabilities fall under external and internal in respect to where they arise relative to the

system. External vulnerabilities are frequently published in news outlets and are growing in

occurrence, while internal ones are far less obvious due to information being held within internal

channels (Baker et al., 2011). Recently, the SolarWinds hack of 2020 has gained a large amount

of attention due to the scale of the attack and the victims. Jibilian & Canales (2021) detail the former by stating that "18,000 of [SolarWinds'] customers installed updates that left them vulnerable to hackers" (para. 6). Some of the victims are Fortune 500 companies such as Microsoft and the attacker is believed to be Russia's Foreign Intelligence Service; thereby showing that even highly capable technology companies are susceptible to external attacks (Jibilian & Canales, 2021). In fact, Baker et al. (2011) attributes 92% of data breaches stemming from external agents, which they state to be a result of an "increase in smaller external attacks rather than a decrease in insider activity" (p.2).
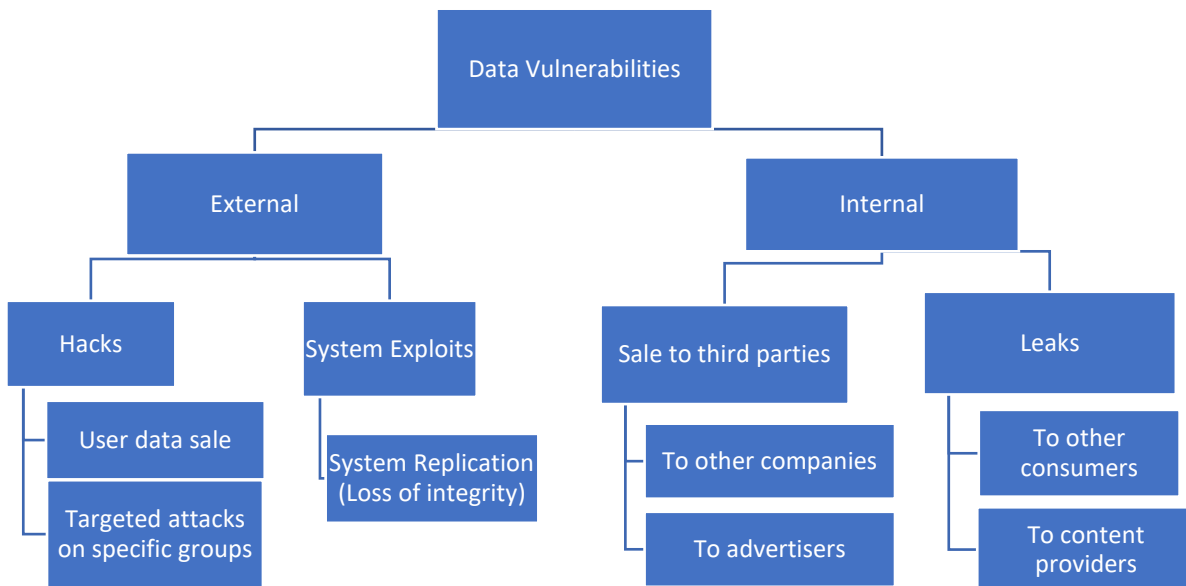


Figure 1. Data Vulnerabilties Within Online Systems: Vulnerabilties can be split into external and internal, where different types of attacks/failures belong to each. Internal vulnerabilities often become a problem due to the information being sent outside of normal channels, while external ones break protections or demonstrate the lack thereof. (Ramshetty, 2021a).

Through the various types of vulnerabilities it is clear that a similarity they share is that user-specific data is prone to being stolen or revealing information about the user as well as the groups they belong to. In the internal case, Figure 1 points out that user data is sold or loses its confidentiality by being moved outside of internal operation. On the other hand with external

vulnerabilities, take for example system exploitation where external agents utilize a known weakness to gain control. In such an attack, there is a possibility that recommendations can be used for reverse engineering of the system, allowing unauthorized actors to mimic and target users within the system (Jeckmans et al., 2013, p. 9). Meaning user data can be used to learn about patterns that were otherwise meant to be protected in order to maintain the privacy of users in respect to the groups they belong to or their preferences.

For solutions of the trade-off to be found it is necessary to locate not just the sources of risks but also those of the vulnerabilities. Others can be observed through some corporations' practice of activities that endanger their users; Jeckmans et al. (2013) identify a few activities in the recommendation system that are a cause for privacy concern, of these activities, "data collection", "data retention", and "recommendations revealing information" stand out (p. 8). To understand why these processes are dangerous, it is crucial to understand what data recommendation systems use and the pipeline that data travels. As Figure 2 depicts, user-specific data such as their ratings on products, preferences/behavior, clicks/searches, etc. influence the recommendation prediction through the model on which the system is built. The data from the user is anlayzed and interpreted by the model according to past inputs and what the model already knows. For each step in the process, data from the user is subject to misuse or mishandling, which is attributed to the activities of data storage and collection that must happen for the system to function. Through data collection the system learns more
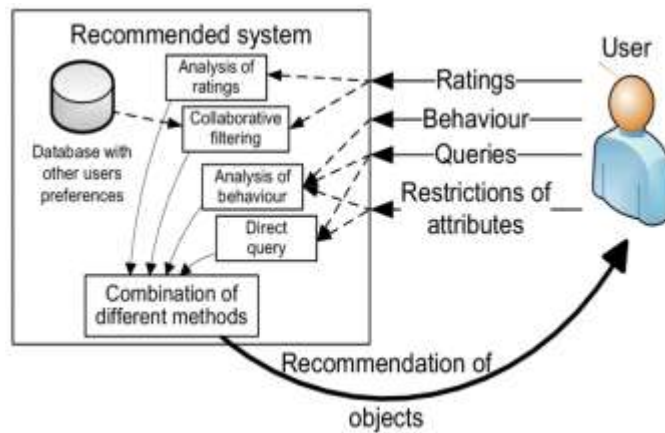


Figure 2. Recommendation Cycle: Flow of data from user to system resulting in a recommendation being output back to user (Eckhardt, 2009, p.61).

about the user and in doing so, creates a connection between the single user and their qualities or wants. On the other hand, data storage maintains a history and catalog of data the system manages. Accordingly, vulnerabilities sprout at the steps in the system where data is utilized, such as at the database shown in Figure 2 on page 5. Therefore, these potential risks to user data security caused by the system's plethora of weaknesses to external or internal agents violate the privacy and confidentiality of consumers.

## CHANGING USER WILLINGNESS TO ACCEPT RISK IN EXCHANGE FOR BETTER RECOMMENDATIONS

Over the countless applications that collect user data in our current technological climate, the fear of the aforementioned data vulnerabilities has not scared away consumers as "the benefits of getting tailored content outweigh any privacy concern" (Harley, 2018, para. 19). Yet, it is unclear whether this sentiment will hold true as recommendation systems become eager to collect more and more data from their users. Firstly, most consumers acknowledge that their data is being stored and "feel various dimensions of control over personal information collection are 'very important' to them" (Madden et al., 2019, para. 4). Thus, given the state of current global consumerism, users have to be cognizant of the services they are using and how each one operates on their data. As a result, Auxier et al. (2019b) report that 81% of Americans feel as though "they have very little/no control over the data [companies] collect," which in turn causes 79% of them to express some concern about the quantity of data collected by these corporations (para. 2; Auxier and Rainie, 2019a). As increasing numbers of users begin to worry over the cost of using recommendations, the likelihood that their utility falls below their needs rises. Most importantly, the trend towards more fear leads to the need for solutions or methods of controlling the data practices employed by current recommendation systems. Popular methods include

regulation of the industry, developing new models such as those in the technical work, and/or consumer consent based data collection.

## REGULATIONS OF THE RECOMMENDATION SYSTEM INDUSTRY MAY BE HARMFUL BUT NECESSARY

The objectives of the recommendation system and that of the public are inherently opposite; the system hopes to provide the best recommendations at the cost of privacy, while the public would rather protect their info and have access to decent recommendations as noted by Harley's (2018) observation that "poor observations are easily ignored … when the benefit of getting good recommendations is strong enough" (para. 29). Figure 3 diagrams Mesthene's technical complexity model, which highlights the relationship between high quality recommendations, low cost of privacy, and wide accessibility of recommendation systems.

Mesthene had used the model to describe how the education system could be widely available and low cost at the expense of quality, which is a common trend in public schools. That is, if public schools were priortized for quality, more money would need to be spent on teachers and resources, which would cause movement upwards and possibly apart from the accessibility and low cost bubbles. Any system
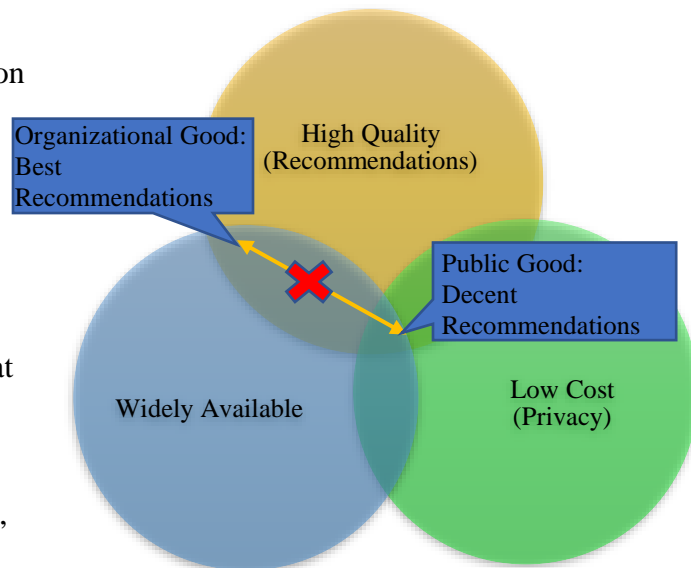
Figure 3. Mesthene's Model - Goals of Recommendations for Public vs. Private: Corporations seeks to improve their recommedaendations due to the proportional increase in profits that follow. The public would rather have all three, which is extremely difficult to meet (Ramshetty, 2021b).

that serves the public or some userbase must reconcile which qualities it values the most, where qualities references the system's accessibiltiy, cost to the user, or quality of product/service. In fact, Mesthene further used healthcare systems to show that they are often expensive since the primary goals are to achieve accessibility and quality. Patients want the best treatment avaiable to them wherever they may be, especially in case of an emergency; therefore, healthcare costs are commonly very expensive due to private firms having to account for the cost of providing the other qualities. The main point being that few systems are able to achieve all three of these qualities and the wants of users versus system managers are vastly different. Some users may prefer lower costs in place of accessibility, but private firms such as those in healthcare may also see the value of profits in not prioritizing cost.

In regard to recommendations, the "x" in the diagram represents the position of the recommendation system in respect to the three categories, where users would like the system to be able to meet each of the three criteria at the center of the diagram. At this location, users receive good recommendations while facing less risk to their personal data because less of their data is used to aid the model in learning preferences. However, the patterns noticed by Adomavicius et al. (2018) demonstrate that the opposite is taking place in today's industry, in which recommendation quality has increased such that users are more biased to the system's output. In other words, system designers have chosen to move towards "organizational good" due to the increase in profits that come with better recommendations, rather than worry about the data they are collecting. Furthermore, the bias Adomavicius et al. (2018) recognizes is the phenomena of recommendations being trusted even if chosen randomly, illuminating why there exists a relationship between profits and recommendation quality. This shift of the "x" to the left in Mesthene's (1970) model depicts the trade-off mentioned earlier in the paper, where the

advancement of a technology comes with some social costs, but there are a multitude of ways to inhibit either extreme of the trade-off. Many see this problem and immediately suggest various regulations or options that the government could use to limit or contain this trend. But, is regulation the best answer to such a problem and should it pertain to the recommendation system alone or all online systems?

Regulations on data have been called for by many agencies since regulation prevents certain activities such as the collection of particular data types; however, the result is often a decrease in recommendation accuracy due to the model's new limited scope of knowledge. An example of regulation harming a system's performance is described by Friedman et al. (2015), they found that a policy in the EU which "severely limited the use of non-essential cookies" caused the online advertising system to be "far less effective in the EU than in other countries" (p. 673). Intuitively, cookies are a necessary component for advertisements because they help track user behavior and preferences, therefore when the policy chooses to protect such sensitive data the system is incapable of providing ads as accurate as before.

Yet, regulators at the government level have pursued legislation that "introduce explicit guidelines and sanctions to regulate data collection, use, and storage," to answer many of the fears that the public expresses (Milano et al., 2020, para. 27). Furthermore, the European Union (EU) has instituted ground-breaking legislation known as General Data Protection Regulation (GDPR), which hopes to deter the "massive collection of personal data about individuals, to the detriment of privacy, but also to a pervasive influence on their behaviour, to the detriment of both individual autonomy and collective interests" (Sartor, 2020, p. 19). However, van Ooijen and Vrabec (2019) find that even such lengthy regulation "fails to solve the problem of information complexity," coming from the advancement of technology as a whole (p. 104).

Hence, it is not possible to continuously resolve tension between better recommendations and privacy/confidentiality through regulation. The research conducted suggests that a shift in the culture of both corporations and society is also needed to achieve a more socially acceptable trade-off balance, one which gives back some control to the users without an abundance of laws.


**USING CONSUMER CONSENT TO MODIFY RECOMMENDATION SYSTEMS**

Another commonly proposed solution is the use of a consent based data collection approach. In this method, systems are keen on asking and involving the user throughout the full process, especially at the stage of data collection. As Arnold (2018) lists, the GDPR requires consent of the user to collect and store data, that consent must be explicit, and users may withdraw consent. Implementations of such a system are varied in style but as Arnold (2018) observes, users that see systems who are "improving their online experiences and offering up information that is relevant to them are more likely to consent to sharing their data in the future" (para. 10).

However, it is unclear how users recognize when systems are improving and whether or not that corporation is instituting the policies it says. Friedman et al. (2015) find three ways in which systems can establish trust: reputation, certification, and trusted computing. Firstly, reputation revolves around the perception of the system over time, since "non-compliance would lead to negative user feedback" (Friedman et al., 2016, 651). Secondly, certification is the same process other industries use to verify compliance to terms, such as certain produce being verified to be organic. Lastly, trusted computing is the idea that the system is able to demonstrate its infrastructure is capable of performing the necessary tasks to uphold its responsibilities (Friedman et al., 2015). Taking these measures allows users to learn to trust the systems they

interact with; if they never do, then they are able to avoid sharing their data to untrusted systems. Nevertheless, these solutions do not answer the question of whether recommendation systems are an outlier in the risk they present to consumers.

## SHOWING THE LARGER SCALE PROBLEM BY COMPARING NOTIFICATIONS AND ONLINE ADVERTISING TO RECOMMENDATIONS

### SIMILARITIES BETWEEN NOTIFICATIONS AND RECOMMENDATIONS

A common online system that is packaged with many applications is the notification system. Ranging from simple messages to life-saving warnings, notifications are what many consider an essential; but, for notifications to be reliable and accurate they must collect data such as contact info, timing information, and priority (Silva et al., 2019, p. 6). These data points are very similar to those of recommendation systems, for example, priority of a notification is related to preferences of a user since they both categorize the user into groups. Additionally, content providers oftentimes intertwine both systems when sending notifications of content suggestions, such as those from Netflix and Amazon.

Referencing back to the terms of privacy and confidentiality, neither system can guarantee a user that their data will remain protected. In respect to privacy, the data users provide may be leaked or hacked and then used to spam certain targeted suggestions utilizing the individual's personal preferences. This same failure of the system also shows how confidentiality of the data has no real structured protection, in neither system protection is a focus or built into the models by design. In other words, the system for recommendations and notifications are not concerned with protecting the data, rather using it for some end-goal. Thereby, it is clear how both systems could be susceptible to the same vulnerabilities and problems.

11

**ONLINE ADVERTISING AND ITS RELATION TO RECOMMENDATIONS**

Another system that plays a large role in the online ecosystem is advertising in the form of website media. The way in which certain ads are targeted to users is through the collection of data including: IP addresses, cookies, web history, etc (Tran, 2014, p. 12-15). In doing so, online systems partake in a data collection activity that resembles that of recommendation systems; both use user data to group individuals together by their preferences. For advertising, web history gives insight into a user's current interests, while for recommendations the user's previous ratings or views would do the same. Though the information collected is not the same, the purpose itself is, which brings to question whether one type of data is more risky than the other. However, in the case of online advertising and recommendations the premise of the data collected is the exact same, both track user activity to predict the next action the user would like to take. Therefore, even if the names of the data are different, if they are used to do the same task then they are directly related and share the same cause for concern.

The vulnerabilities mentioned throughout the paper also extend to advertising; of the sites studied in Tran's (2014) paper, "55% … directly leak [a] piece of private information" (p. 29). Moreover, regardless of external or internal agents, notifications, advertising, and recommendations all face the same problems due to the nature of the data they store. For example, if a user visits Amazon and searches through a variety of shoes, both advertising and Amazon's recommendation engine will show the user further shoe options. Thereby, the data leaked/hacked from either system will depict the same preferences of that individual. Consequently, a system that manages data which is valuable to third-parties must build a purposeful protection infrastructure that works with the social solutions of consent based data collection and regulation detailed earlier in the paper. Hence, not only are these two systems just

as vulnerable as recommendations, hundreds of other online systems that collect data without addressing the storage and protection of said data are sources of risk.

**IMPORTANCE OF DATA PROTECTIONS IN RECOMMENDATION SYSTEMS AND BEYOND**

The future of the recommendation system market depends on the ability of the engineers to find new advancements of the science and for users to not feel at risk when using the system. Recommendations are not more risky than systems such as notifications and online advertising, due in part to the shared data activities and lack of importance to how data is handled. Furthermore, with the information that regulation's efficacy is questionable, studies in the future may find a different approach to alleviating the trade-off between advancing technology and social good. One method may be to validate a data trait, know as the "right to be forgotten"; within such a system, "recommender providers are obliged to delete the personal data any time the when the data subject requests it" (Tejada-Lorente et al., 2018, p. 6). With the ability for a user to control their data's existence in databases, it is expected that they will be more trusting with such systems and can be compared to the outcomes of asking them for consent.

Another manner to promote data security could be to require an enforcing agency that is part of each system or that monitors an array of systems (Tejada-Lorente et al., 2018, p. 6). The goal of the monitoring entity would be to enforce data policies that are made at higher levels or agreed upon within an industry. Through this method, any violations or breaches would be subject to immediate review, followed by required improvements and/or penalties. Without proactive measures against the current trend away from consumer data safety, many of today's online systems are susceptible to sharp consequences resulting from their inaction. Some of these consequences could be the loss of users and fines from the government, which would negatively

impact the entire online ecosystem. Thus, there is an obvious need for governments,

corporations, and society to work together to define proper manners in which data is managed

among the three. All in all, data will continue to influence every facet of our lives and

recommendation systems are just a small piece of an intricate everchanging online machine.

# WORKS CITED

Adomavicius, G., Bockstedt, J., Curley, S. P., Zhang, J., & Ransbotham, S. (2018, November 13). The hidden side effects of recommendation systems. *MIT Sloan Management Review.* Retrieved from https://sloanreview.mit.edu/

Arnold, A. (2018, May 7). How GDPR will impact content recommendation engines. *Forbes.* Retrieved from https://www.forbes.com/

Auxier, B., & Rainie, L. (2019a, November 15). Key takeaways on Americans' views about privacy, surveillance and data-sharing. *Pew Research Center.* Retrieved from https://www.pewresearch.org/

Auxier, B., Rainie, L., Anderson, M., Perrin, A., Kumar, M., & Turner, E. (2019b, November 15). Americans and privacy: concerned, confused, and feeling lack of control over their personal information. *Pew Research Center.* Retrieved from https://www.pewresearch.org/

Baker, W., Hutton, A., Hylender, D. C., Pamula, J., Porter, C., Spitler, M., … Valentine, A. J. (2011). 2011 Data breach investigations report. *Verizon*. Retrieved from https://cybersecurity.idaho.gov/wp-content/uploads/sites/87/2019/04/data-breach-investigations-report-2011.pdf

Bennett, J., & Lanning, S. (2007) The Netflix prize. (n.p.). Retrieved from https://www.cs.uic.edu/~liub/KDD-cup-2007/NetflixPrize-description.pdf

Eckhardt, A. (2009). Various aspects of user preference learning and recommender systems. [1]. *Proceedings of the Dateso 2009 Annual International Workshop on Databases, 471*, 56-67. Retrieved from https://www.semanticscholar.org/

Faggella, D. (2017, October 30). The ROI of recommendation engines for marketing. *Martech.* Retrieved from https://martechtoday.com/

Friedman, A., Knijnenburg, B., Vanchecke, K., Martens, L., & Berkovksy, S. (2015). Privacy aspects of recommender systems. *Recommender Systems Handbook,* 649-688. https://doi.org/10.1007/978-1-4899-7637-6_19

Harley, A. (2018, September 30). Individualized recommendations: users' expectations & assumptions. Nielsen Norman Group. Retrieved from https://www.nngroup.com/

Jeckmans, A. J. P., Beye, M., Erkin, Z., Hartel, P., Lagendijk, R. L., & Tang, Q. (2013). Privacy in recommender systems. *Social Media Retrieval,* 263-281. doi:10.1007/978-1-4471-4555-4_12

Jibilian, I., & Canales, K. (2021, February 25). Here's a simple explanation of how the massive SolarWinds hack happened and why it's such a big deal. *Business Insider.* Retrieved from https://www.businessinsider.com/

Madden, M., & Rainie, L. (2019, December 31). Americans' views about data collection and security. *Pew Research Center*. Retrieved from https://www.pewresearch.org/

Martin, M. W., & Schinzinger, R. (2010). *Introduction to engineering ethics* (2nd ed.). New York, NY: McGraw-Hill.

Mesthene, E. G. (1970). *Technological change: its impact on man and society*, 63-89. Cambridge, MA: Harvard University Press.

Milano, S., Taddeo, M., & Floridi L. (2020, February 27). Recommender systems and their challenges. *AI & Soc*, *35*, 957–967. Retrieved October 06, 2020, from https://link.springer.com/

Ramshetty, S. (2021a). The variety of data vulnerabilities within online systems. [2]. *STS Research Paper: Comparison and evaluation of user privacy risk in recommendation systems* (Unpublished undergraduate thesis). School of Engineering and Applied Science, University of Virginia. Charlottesville, VA

Ramshetty, S. (2021b). Mesthene's model - goals of recommendations for public vs. private. [3]. *STS Research Paper: Comparison and evaluation of user privacy risk in recommendation systems* (Unpublished undergraduate thesis). School of Engineering and Applied Science, University of Virginia. Charlottesville, VA

Sartor, Giovanni (2020, June). The impact of the general data protection regulation (GDPR) on artificial intelligence [PDF File]. *Scientific Foresight Unit*, 15-79. Retrieved from https://www.europarl.europa.eu/RegData/etudes/STUD/2020/641530/EPRS_STU(2020)641530_EN.pdf

Silva, L. A., Leithardt, V. R. Q., Rolim, C. O., Villarrubia, G., Geyer, C. F. R., & Silva, J. S. (2019). PRISER: Managing notification in multiple devices with data privacy support. *Sensors 2019*, *19(14)*. doi:10.3390/s19143098

Tejada-Lorente, A., Bernabe-Moreno, J., Herce-Zelaya, J., Porcel, C., & Herrera-Viedma, E. (2018). Adapting recommender systems to the new data privacy regulations. (n. p.). doi: 10.3233/978-1-61499-900-3-373

Tran, M. (2014). Privacy challenges in outline targeted advertising. *Computers and Society*. Retrieved from https://tel.archives-ouvertes.fr/tel-01555362/document

van Ooijen, I., & Vrabec, H. U. (2019). Does the GDPR enhance consumer's control over personal data? An analysis from a behavioural perspective. *J Consum Policy, 42*, 91–107. doi:10.1007/s10603-018-9399-7