

# PRIVACY PRESERVING TECHNIQUES AND THEIR APPLICATIONS IN ELEARNING

***Malinka Ivanova, Iskra Trifonova***

*Technical University of Sofia, Faculty of Applied Mathematics and Informatics,  
Department of Informatics, Bulgaria*

[mivanova@tu-sofia.bg](mailto:mivanova@tu-sofia.bg); [izi1976@abv.bg](mailto:izi1976@abv.bg)

## ТЕХНИКИ ЗА ЗАПАЗВАНЕ НА ПОВЕРИТЕЛНОСТТА И ТЕХНИ ПРИЛОЖЕНИЯ В ЕЛЕКТРОННОТО ОБУЧЕНИЕ

**Abstract:** *The paper summarizes contemporary methods and techniques for privacy preservation as some challenging issues are analyzed and presented. A bibliometric approach is utilized in order for the “big picture” to be outlined, showing current research status and trending topics. The bibliographic data are taken from scientific database Scopus and processed through specialized software. In addition, a detailed review is also performed to classify problems and solutions in the area of privacy preservation. Special attention is given to possibilities for data privacy protection in intelligent eLearning environments. The role of machine learning for creating more secure data models is pointed out. A conceptual model, summarizing the findings, is proposed.*

**Keywords:** *Privacy Preservation; Machine Learning; eLearning Intelligent Environment*

### Въведение

Понастоящем се събира огромно количество данни чрез уеб или мобилни приложения, които впоследствие се обработват и анализират за различни цели. В голяма част от случаите тези данни са лични и чувствителни и използването им от атакуващ би нанесло вреди на потребителите или би нарушило тяхното лично пространство. Поради тази причина, отделни изследователи и изследователски екипи търсят подходи за предпазване и защита на поверителността. Добре би било, преди данните да бъдат изпратени за анализ, да бъдат модифицирани при спазване на баланс между запазване на поверителността на лицата и получаване на висока степен на използваемост на данните. Една част от съществуващите алгоритми, методи и

техники са само теоретични разработки, докато друга част успешно се прилагат на практика. Въпреки, постигнатото до момента в областта на защита на поверителността, академията и индустрията продължават да изследват тази област, за да се идентифицират добри подходи с практическа приложимост.

Съвременните системи за електронно обучение се характеризират с елементи на изкуствен интелект, улеснявайки преподавателя и подпомагайки студентите в тяхното обучение. В тази връзка се събират различни данни: лични, чувствителни, касаещи предишно образование и учебен успех, както и текущи резултати, като много често студентът не е информиран по какъв начин и за какви цели се използват те. Не са предвидени механизми за избор, чрез които студентът да прецени коя информация да бъде споделена и каква част да бъде скрита. Освен това, някои данни (например за успеваемостта на студентите) могат да се използват извън системите за електронно обучение и получената статистика да стане публично достъпна. Това улеснява атакуващ, който при наличие на допълнителни информационни източници, може да проведе различни по вид атаки за откриване идентичността на конкретно лице или на чувствителни данни за него.

Целта на статията е да направи преглед и анализ на съвременното положение, като обобщава научни постижения и добри практики за запазване и защита на поверителността. Създаден е концептуален модел, представящ актуални подходи, методи и техники за защита на лични и чувствителни данни.

Най-напред е извършен библиометричен анализ за очертаване на глобалната картина и посочване на интересните за момента теми за изследване. След това, е извършен анализ на съдържанието на научни статии за по-добро разбиране на постиженията и получените резултати от научната общност. Представен е концептуален модел, обобщаващ библиометричния поглед и литературния анализ. Накрая са направени изводи.

## Библиометричен анализ

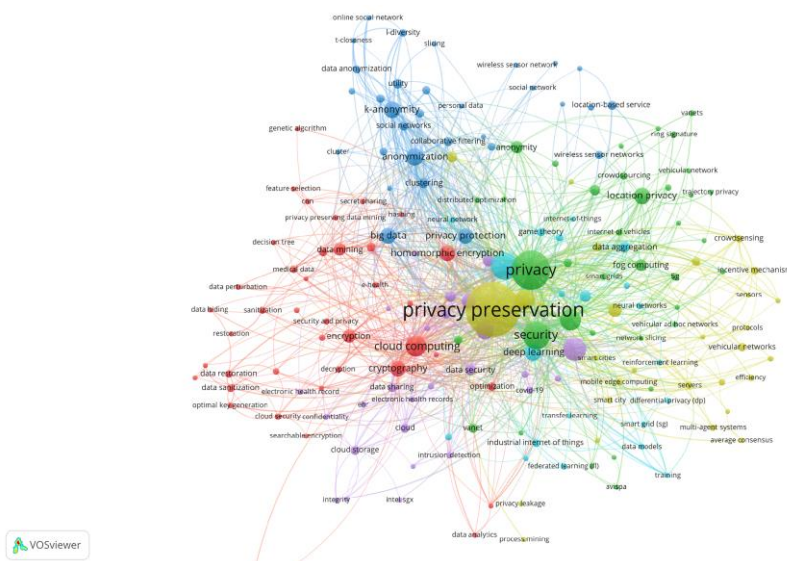
За извършване на библиометричен анализ е подадена заявка в научната база данни Scopus на 25 март 2023 г., състояща се от израза: *privacy preservation*. Търсенето е извършено в заглавието, абстракта и ключовите думи на индексирани статии, глави от книги и книги. Върнатият резултат посочва 5447 документа, отговарящи на заявката. Търсенето е стеснено чрез изследване на документите, които са публикувани през последните пет години и тези, индексирани през тази година, т.е. от 2018 г. до настоящия момент, както и се интересуваме от документи, написани на английски език. Полученият резултат след критериалното стеснение включва 3331 документа. Допълнително е извършено сортиране на документите по уместност. Библиографската информация за 2000 от най-уместните документи е използвана за по-нататъшно проучване чрез VOSviewer [7].

Най-напред, за извършване на анализ на включените от авторите на документите ключови думи, е създадено визуално представяне на съвместна мрежа. Съвместната мрежа показва използваните по двойки термини от ключовите думи,

като тя е изградена при зададен минимален брой на срещане на тези ключови думи, който да бъде 5. Това означава, че поне в 5 документа се среща съвместно една двойка ключови думи. Съвместната мрежа показва и графично връзката между тези два термина. При конструирането ѝ се изчислява *честотата* на съвместно срещаните термини и се намират централни термини, около които се формират тематични клъстери. Друг параметър на съвместната мрежа е *силата на връзката*, който показва броя съвместни срещания на един термин с други термини.

За подадената заявка *privacy preservation*, конструираната мрежа на съвместно използваните термини е показана на Фигура 1. Формирани са шест тематични клъстера с централни термини: *privacy preservation*, *privacy*, *federated learning*, *anonymization*, *differential privacy* и *cloud computing*.

1. По-често срещаните термини в клъстера *privacy preservation* са: *Internet of things*, *edge computing*, *data publishing*, *multi-agent systems*, *smart city*, *sensors*, *protocols*, други.
2. Клъстерът *privacy* включва термини като *security*, *authentication*, *location privacy*, *anonymity*, *fog computing*, *GDPR* и други.
3. Около термина *federated learning* е образуван друг клъстер с ключови думи като: *healthcare*, *electronic health records*, *smart contract*, *machine learning*, *artificial intelligence*, *access control* и други.
4. Клъстерът с централна дума *anonymization* е образуван около термините *big data*, *k-anonymity*, *l-diversity*, *t-closeness*, *generalization*, *information loss*, *clustering* и други.
5. Централният термин *differential privacy* формира клъстер с ключовите думи: *deep learning*, *aggregation data*, *smart grid*, *internet of vehicles*, *reinforcement learning*, *face recognition* и други.
6. Клъстерът *cloud computing* се състои от термините: *cryptography*, *homomorphic encryption*, *data perturbation*, *key generation*, *feature selection*, *secret sharing*, *association rule mining*, други.



Фигура 1. Конструирана съвместна мрежа при подадена заявка *privacy preservation*

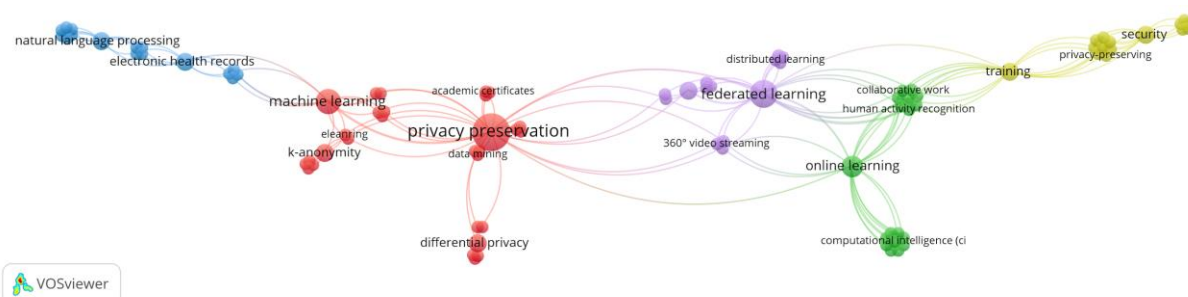
Първите 24 най-често съвместно срещани ключови думи и силата на връзката са показани в Таблица 1. Почернените термини са централни и около тях са формирани клъстери.

Термин	Честота	Сила на връзката	Термин	Честота	Сила на връзката
<b>privacy preservation</b>	560	974	cloud computing	81	177
<b>privacy</b>	290	656	internet of things	77	171
<b>security</b>	143	424	machine learning	59	158
<b>differential privacy</b>	143	253	cryptography	47	146
<b>authentication</b>	87	220	anonymization	59	137
<b>federated learning</b>	101	216	deep learning	52	124

Таблица 1. Честота и сила на връзката на най-често срещаните ключови думи

При подадена заявка *privacy preservation and e-learning* за периода от 2018 г. до момента като резултат са върнати само 24 документа. Това показва, че тази тема е в ранен стадий на изучаване и са необходими бъдещи изследвания. Създадената съвместна мрежа е показана на Фигура 2 и в нея са оформени 5 клъстера:

1. Най-големият клъстер е образуван около термина *privacy preservation* и включва термините: *e-learning, mobile-based learning system, distributed online learning, differential privacy, k-anonymity, anonymization, digital footprints, machine learning, data mining*.
2. Централна дума за втория клъстер е *federated learning*, който се състои от термините: *distributed learning, edge computing, virtual reality, 360° video streaming, digital twin, internet of things*.
3. Третият клъстер е формиран около термина *online learning* и е свързан с термините: *computational intelligence, collaborative work, human activity recognition, privacy, real-time systems, semi-supervised learning, data models*.
4. Термините *training* и *security* са с еднаква честота в четвъртия клъстер, който включва термините: *privacy-preserving, transfer learning, distributed databases, online teaching*.
5. В последният клъстер термините *natural language processing* и *electronic health records* са с еднаква честота на срещане в документите, като останалите термини в този клъстер са: *privacy protection, confidentiality, generative language models, ensemble classification, heuristic-based perturbation*.



Фигура 2. Конструирана съвместна мрежа при подадена заявка privacy preservation and e-learning

Термин	Честота	Сила на връзката	Термин	Честота	Сила на връзката
privacy preservation	9	32	natural language processing	2	11
federated learning	5	27	electronic health records	2	10
machine learning	4	16	digital twin	2	8
online learning	3	24	k-anonymity	2	8
training	2	19	differential privacy	2	6
security	2	14	collaborative work	1	10

Таблица 2. Честота и сила на връзката на най-често срещаните термини при заявка privacy preservation and e-learning

## Анализ на научни публикации

След библиометричния анализ е извършено и проучване на съдържанието на наскоро публикувани научни документи, за получаване на по-добро разбиране на съвременната проблематика.

Оптимизиран алгоритъм за запазване на поверителността, основан на метода k-анонимност, е предложен от Mahapatra и др. [6]. Алгоритъмът е разработен за прилагане върху специален вид набори от данни, при които всички квази-идентификатори се намират в един и същ домейн. Наборите с данни се характеризират с идентична генерализационна йерархия, за разлика от съществуващите алгоритми, които работят с общи такива. Алгоритъмът е с подобро евристично търсене и бързо намира оптималното решение на k-анонимност. Gangarde и др. разработват модел за защита на чувствителни данни, събирани в онлайн социални мрежи [3]. Той включва три усъвършенствани фази: за клъстериране на данните, за постигане на k-анонимност и за прилагане на техниките l-разнообразие и t-близост. Авторите показват, че предложеният модел е по-добър в сравнение със съществуващи решения, като параметрите: степен на анонимизация, загуба на информация и изчислителна ефективност са подобрени. Предложеното

решение има и някои ограничения, като например, че работи добре само със статични набори от данни.

Защитата на поверителността на мобилни потребители се дискутира от Sun и др., като е разработен алгоритъм за генериране на фиктивни местоположения, които изключват данни за оригиналната локация [9]. Алгоритъмът е основан на  $k$ -анонимност и  $1/2$  разнообразие и включва изпълнението на две процедури: (1) избор на  $k$ - $n$  фиктивни местоположения въз основа на техните стойности на параметъра ентропия и постигане на  $k$ -анонимност и (2) разглеждане на сценарий, в който избрани  $k$  местоположения включват оригиналната локация. Демонстрирана е устойчивост срещу (1) колизионна атака, при която атакуващ влиза в сътрудничество с други потребители или доставчик на приложения за услуги, базирани на местоположение и разкрива търсеното местоположение и (2) дедуктивна атака, която се предполага, че се извършва от доставчик на приложения за услуги, базирани на местоположение, който има познание за предишни и настоящи заявки на потребителя, както и познава защитната схема.

Едно решение за защита на чувствителни данни, които се събират чрез устройства, предназначени за Интернет на нещата и в контекста на „умен град“, е представено от Gheissari и др. [4]. За защита на поверителността е използван методът диференциална поверителност с вероятно разпределение на Лаплас или експоненциално разпределение и е приложена парадигмата софтуерно дефинирана мрежа. Авторите считат, че предложеното от тях решение по-ефективно защитава срещу разкриване на информация в сравнение с приложен метод диференциална поверителност върху статични данни. Truong и др. обобщават и дискутират съвременни подходи за запазване на поверителността, като анонимизиране на данните, диференциална поверителност, протокол за многостранно изчисление и хомоморфно криптиране [12]. Особено внимание обръщат на техниката федеративно обучение и предизвикателни въпроси при прилагане на Общия регламент относно защитата на данните, когато се създават модели чрез машинно обучение. Обобщени са и атаките срещу техниките за запазване на поверителността при федеративното обучение.

Преглед на разработени подходи за защита на чувствителни данни при използване на устройства за Интернет на нещата е направен от Torge и др. [11]. Подходите са категоризирани в осем групи: техники за анонимизиране (включително  $k$ -анонимност,  $l$ -разнообразие,  $t$ -близост), за объркване, машинно обучение на няколко нива, децентрализирано машинно обучение, хомоморфно криптиране, създаване на модели на поток от данни с различни нива на поверителност, създаване на сбита версия на данните и създаване на сигурни индивидуални хранилища на данни. Една част от подходите са само теоретични и все още не се прилагат на практика. Rafiei и van der Aalst дискутират групово-базирани техники за запазване на поверителността при извличане на информация от лог файлове [8]. Предложен е разширен вариант на метода *TLKS*-поверителност, който използва техники като: точност на времевия отпечатък ( $T$ ), дължина на изречението ( $L$  – значение на основното познание),  $k$ -анонимност ( $K$ ) и граница на доверието относно



чувствителните атрибути (С). Експериментите показват, че разширената версия на TLKS-поверителност запазва по-добре поверителността от предишната. Въпреки това, в лог файлове на събития с голямо съотношение на уникални следи и когато предполагаемият тип основни познания е много специфичен, груповите техники за защита на поверителността може да не са в състояние да запазят общата полезност на данните.

Един подход за постигане на по-висока степен на поверителност на данните при предаване в автомобилна ад хок мрежа е представен в [1]. Авторите предлагат сигурна схема за автентикация, използваща псевдоними. Схемата премахва ролята на крайпътни единици, пред които автомобилите да се автентикират и предлага директна автомобилна комуникация. Установено е, че подходът е устойчив на вътрешни атаки, тъй като доверителна единица може да отмени сертификатите на измамни превозни средства, като им попречи да излъчват непрекъснато фалшиви съобщения. Схемата използва концепции, заложи в криптография, основана на елиптични криви и хеш криптографски функции. Като бъдеща работа е посочено проектиране на автентикационна схема без използване на криптография, основана на елиптични криви в 5G автомобилна мрежа, а по-скоро ще бъдат използвани техники от изчисления в мъгла.

Теки и др. предлагат модифициран подход за запазване на поверителността, в който се прилагат алгоритми от машинното обучение за построяване на регресионни и класификационни дървета (Classification and Regression Tree (CART)) и нереализиран алгоритъм [10]. Стойностите на атрибутите в базата данни се изкривяват и се премахват дублиращите записи, като се получават набор с разбъркани записи и нереализиран набор с данни. Върху оригиналните данни се прилага алгоритъмът CART, а върху новополучените два набора модифициран CART (MCART). Ако при сравнение на построените дървета се получи един и същ резултат на стойностите Gini, то е постигната защитеност на данните. В бъдеще този алгоритъм може да бъде подобрен чрез прилагане на усъвършенствани техники за класификация като C4.5, CHAID, SVM и др. Ezhilarasi и др. разработват генератор на синтетични данни, наречен FBprophet, с възможности за прогнозиране на времеви серии с данни [2]. Генерираните синтетични данни са подходящи за последваща обработка чрез алгоритми от машинното обучение. Така получените модели се характеризират с висока степен на защитеност на поверителността. Изследване и анализ за възможни приложения на техники за запазване на поверителността в електронното обучение е представено в [5]. Направени са експерименти върху набор с данни от реален учебен процес чрез алгоритми за k-анонимност и  $(\epsilon, \delta)$ -диференциална поверителност. Машинно обучение е използвано за прогнозиране на подходящ модел за запазване на поверителността.

## Концептуален модел

Обобщение на получените резултати от библиометричния анализ и подробното проучване на съдържанието на научни статии е представено чрез създаден концептуален модел (Фигура 3). Най-често дискутираните подходи за защита на поверителността са класифицирани в пет групи: (1) основани на анонимизация и

псевдо анонимизация, (2) използващи криптографски методи, (3) генериране на синтетични данни, (4) методи за защита на модели, създадени чрез машинно обучение и (5) други техники. Търсят се подходи за защита както на статични бази данни, така и на динамично пристигащи потоци от данни. Сред приложенията на тези подходи особен интерес се проявява към защита на лични и чувствителни данни във връзка с изграждане на умни градове, Интернет на нещата и използване на сензорни устройства, в автономните транспортни средства, здравеопазване, бизнес и търговия, електронно правителство, мобилни комуникации. Напоследък, се дискутират и възможности за използване на определени подходи в среди за електронно обучение за предпазване от разкриване на данни за участници в образователен процес.



Фигура 3. Концептуален модел относно съвременни подходи за запазване на поверителността



## Заклучение

Извършеното библиометрично и литературно проучване в статията показва силна необходимост от разработване на подходи за запазване на поверителността, което произтича от непрекъснатото събиране на лични и чувствителни данни в уеб пространството, от мобилни приложения, в сензорни и други видове автономни мрежи. Предложени са различни методи и техники, като сред най-често дискутираните са: техники за анонимизиране и псевдо анонимизиране, криптографски методи, използване на синтетични данни. Специално внимание се отделя на защитата на създадени чрез машинно обучение модели, като най-популярни подходи са: федеративното обучение, децентрализирано машинно обучение и машинно обучение на няколко нива. В електронното обучение разработките относно прилагане на подходи за защита на поверителността са в начален стадий и са необходими по-нататъшни изследвания.

## Acknowledgments

This research is supported by the Bulgarian FNI fund through the project "Modeling and Research of Intelligent Educational Systems and Sensor Networks (ISOSeM)", contract КП-06-Н47/4 from 26.11.2020.

## References // Литература

- [1] Al-Shareeda, M. A.; Anbar, M.; Manickam, S.; Hasbullah, I. H. (2022). "A Secure Pseudonym-Based Conditional Privacy-Preservation Authentication Scheme in Vehicular Ad Hoc Networks", *Sensors (Basel)*, 22(5):1696, 2022. DOI: <https://doi.org/10.3390/s22051696>
- [2] Ezhilarasi, P.; Ramesh, L.; Xiufeng, L.; Jens Bo, H.-N. (2023). "Smart Meter Synthetic Data Generator development in python using FBProphet", *Software Impacts*, vol.15, 2023. DOI: <https://doi.org/10.1016/j.simpa.2023.100468>
- [3] Gangarde, R.; Sharma, A.; and Pawarq, A. (2023). "Enhanced Clustering Based OSN Privacy Preservation to Ensure k-Anonymity, t-Closeness, l-Diversity, and Balanced Privacy Utility", *Computers, Materials & Continua*, 75(1), pp. 2171-2190, 2023. DOI: <https://doi.org/10.32604/cmc.2023.035559>
- [4] Gheisari, M.; Shojaeian, E.; Javadpour, A.; Jalili, A.; Esmaeili-Najafabadi, H.; Bigham, B. S.; Vorobeva, A. A.; Liu, Y.; Rezaei, M. (2023). "An Agile Privacy-Preservation Solution for IoT-Based Smart City Using Different Distributions", in *IEEE Open Journal of Vehicular Technology*, vol. 4, pp. 356-362, 2023, DOI: <https://doi.org/10.1109/OJVT.2023.3243226>
- [5] Ivanova, M.; Trifonova, I.; and Bogdanova, G. (2022). "Privacy Preservation in eLearning: Exploration and Analysis", 2022 20th International Conference on

- Information Technology Based Higher Education and Training (ITHET), Antalya, Turkey, 2022, pp. 1-8, 2022. DOI: <https://doi.org/10.1109/ITHET56107.2022.10031904>
- [6] Mahanan, W.; Chaovalitwongse, W. A.; and Natwichai, J. (2021). “Data privacy preservation algorithm with k-anonymity”, World Wide Web 24, pp. 1551–1561, 2021. DOI: <https://doi.org/10.1007/s11280-021-00922-2>
- [7] Perianes-Rodriguez, A.; Waltman, L.; and Van Eck, N. J. (2016). “Constructing bibliometric networks: A comparison between full and fractional counting”. Journal of Informetrics, 10(4), pp. 1178-1195, 2016. DOI: <https://doi.org/10.1016/j.joi.2016.10.006>
- [8] Rafiei, M.; and van der Aalst, W. M. P. (2021). “Group-based privacy preservation techniques for process mining”, Data & Knowledge Engineering, vol. 134, 2021. DOI: <https://doi.org/10.1016/j.datak.2021.101908>
- [9] Sun, G.; Cai, S.; Yu, H.; Maharjan, S.; Chang, V.; Du, X.; and Guizani, M. (2019). “Location Privacy Preservation for Mobile Users in Location-Based Services”, IEEE Access, vol. 7, pp. 87425-87438, 2019, DOI: <https://doi.org/10.1109/ACCESS.2019.2925571>
- [10] Teki, S. M.; Banothu, B.; Varma, M. K. (2019). “An Un-realized Algorithm for Effective Privacy Preservation Using Classification and Regression Trees”, Revue d'Intelligence Artificielle, Vol. 33, No. 4, pp. 313-319, 2019. DOI: <https://doi.org/10.18280/ria.330408>
- [11] Torre, D.; Chennamaneni, A.; and Rodriguez, A. (2023). “Privacy-Preservation Techniques for IoT Devices: A Systematic Mapping Study”, IEEE Access, vol. 11, 16323-16345, 2023. DOI: <https://doi.org/10.1109/ACCESS.2023.3245524>
- [12] Truong, N.; Sun, K.; Wang, S.; Guitton, F.; and Guo, Y. (2021). “Privacy preservation in federated learning: An insightful survey from the GDPR perspective”, Computers & Security, vol. 110, 2021. DOI: <https://doi.org/10.1016/j.cose.2021.102402>

Received: 24-04-2023

Accepted: 29-06-2023

Published: 24-07-2023

Cite as:

Ivanova, M.; Trifonova, I. (2023). “Privacy Preserving Techniques and Their Applications in Elearning”, Science Series “Innovative STEM Education”, volume 05, ISSN: 2683-1333, pp. 93-102, DOI: <https://doi.org/10.55630/STEM.2023.0512>