

▣ Problematika ohranjanja zasebnosti pri podatkovnem rudarjenju dokumentov z občutljivimi podatki

Matjaž Kragelj, National and University Library, Turjaška 1, 1000 Ljubljana, Slovenia matjaz.kragelj@nuk.uni-lj.si
 Mirjana Kljajić Borštnar, University of Maribor, Faculty of Organizational Sciences, Kranj, Slovenia, e-mail: Mirjana.Kljajic@um.si
 Alenka Brezavsček, University of Maribor, Faculty of Organizational Sciences, Kranj, Slovenia, e-mail: Alenka.Brezavscek@um.si

Izvleček

V prispevku obravnavamo problem, s katerim se soočamo pri uporabi dokumentov, ki poleg vsebinskih podatkov vsebujejo tudi občutljive podatke o posamezniku, ki omogočajo njegovo razkritje tudi, ko to ni zaželeno. Med področja, kjer nastane veliko podatkov te vrste, štejeemo zdravstveno varstvo, transport, kazenski pregon in nacionalno varnost, izobraževanje, sodobne internetne storitve, področje sodobnih aplikacijskih ekosistemov, internet stvari, finančni sektor in odprte podatke državne uprave. Cilj je zaščititi zasebnost subjekta ter hkrati zagotoviti kakovostne podatke za nadaljnje poglobljene analize in s tem nudenje novih znanj za naprej. Za reševanje omenjenih izzivov na področju podatkovnega rudarjenja se je razvilo posebno podpodročje, imenovano PPDM – Privacy Preserving Data Mining, ki se ukvarja z ohranjanjem zasebnosti pri tem procesu. Sistematično smo pregledali relevantno literaturo podpodročja PPDM in opisali glavne metode in tehnike. Tehnike PPDM so zasnovane tako, da zagotavljajo določeno raven zasebnosti, obenem pa ohranjajo uporabnost podatkov, da se lahko uporaba še vedno učinkovito izvaja na transformiranih podatkih. Metode, s katerimi dosegamo zaščito posameznika na eni in uporabno vrednost podatkov na drugi strani v grobem delimo na metode razprševanja podatkov, metode izkrivljanja (z uporabo anonimizacije, randomizacije, vrtenja in vnašanjem šuma v podatke) ter metode šifriranja podatkov. Za doseganje višje zaščite lahko uporabimo tudi kombinacije teh metod. Poleg pregleda metod smo podali nekaj praktičnih primerov ter našli domene oz. področja, kjer se kaže potreba po nadaljnji analizi in ponovni uporabi podatkov, a hkrati potreba po anonimizaciji oz. prikritju lastnika (subjekta) in njegovih podatkov (atributov).

Ključne besede: podatkovno rudarjenje, osebni podatki, ohranjanje zasebnosti, metode ohranjanja zasebnosti podatkov, pregled literature, varnost podatkov

1. UVOD

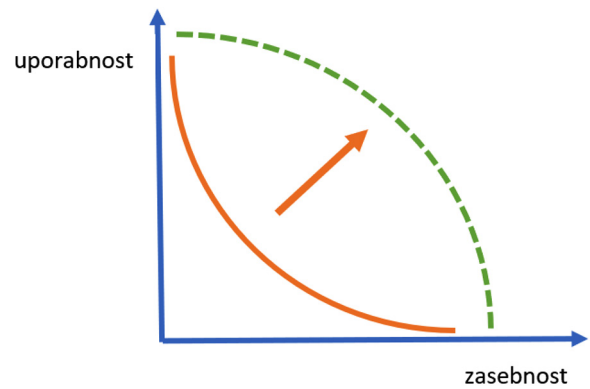
S pojavom interneta se je pojavila potreba in možnost po orodjih za iskanje, razvrščanje in kasneje analizo zbranih digitalnih podatkov. Kot primer lahko navedemo, da sta še v prejšnjem stoletju Smith & Chang (1997), razvila spletno orodje – WebSeek, ki je (bilo) namenjeno iskanju in sortiranju slik s spleta. Dokler so bile količine podatkov obvladljive, so bile tudi metode in orodja razmeroma enostavna, v današnjem času pa potreba po zbiranju in obdelavi podatkov strmo narašča [2]. V podjetju IBM ugotavljajo, da je bilo v letih 2012 in 2013 ustvarjeno več kot 90 % vseh podatkov na

svetu [3], do leta 2025 pa bo za analizo primernih več kot 150 zetabaj-tov (10^9 terabajtov). Z naraščanjem količine digitalnih podatkov, ki jih je potrebno obdelati, se je pojavila potreba po vpeljavi naprednejših metod, ki temeljijo

na principih umetne inteligence in strojnega učenja, natančneje – podatkovnega rudarjenja. Gre za proces pridobivanja implicitnih informacij in znanj, ki bi lahko bile koristne, črpanje le-teh pa poteka iz množičnih, neurejenih, nepopolnih, nejasnih ali naključnih, ne nujno strukturiranih podatkovnih struktur (Sahu, Shrima & Gondhalakar, 2008).

Med digitalnimi in digitaliziranimi dokumenti, ki jih rudarimo, so tudi takšni, ki so občutljive narave in zaradi tega zahtevajo posebno skrb in previdnost pri hranjenju, obdelavi in posredovanju tretjim osebam. Neustrezno postopanje pri rudarjenju takih dokumentov lahko povzroči grožnjo zasebnosti preko razkritja identitete in s tem posredno povezanih podatkov (atributov), saj attribute navadno povezujemo z lastništvom ravno preko identitete [5]. V skladu z definicijo iz slovarja Cambridge, je zasebnost definirana kot pravica posameznika, da obdrži svoje osebne podatke, zadeve in odnose v tajnosti [6]. Definicija sledi 12. členu Splošne deklaracije o človekovih pravicah, ki pravi, da se ni dovoljeno nikomur samovoljno vmešavati v zasebno življenje, družino, dom ali dopisovanje, prav tako pa ni dovoljeno žaljenje časti in dobrega imena slehernega posameznika. Vsakdo ima pravico do pravnega varstva pred takšnim vmešavanjem ali napadi [7]. Težava, ki jo pri rudarjenju dokumentov z občutljivo vsebino zaznavamo, je naslednja: analizo želimo izvajati tako, da v popolnosti ohranjamo vrednosti in pomen podatkov, ki pomenijo attribute, a hkrati želimo zaščititi občutljive podatke o posamezniku, torej zaščititi njegovo zasebnost. Kot ugotavljajo Gokulnath et al. (2015), je ohranjanje zasebnosti pri rudarjenju občutljivih in osebnih podatkov ključnega pomena za učinkovito izvedeno podatkovno rudarjenje.

V drugem poglavju bomo navedli glavna področja, kjer se ustvarjajo podatki in dokumenti z občutljivimi atributi. Med temi področji veljajo podatki o zdravstvenem varstvu za najpomembnejše, a so hkrati najbolj občutljivi, saj vsebujejo vse zasebne podatke, ki so informacije o pacientu, kot so bolezni, podatke o zdravljenju, recept, ime, naslov itd. Takšna zbirka podatkov katerekoli zdravstvene organizacije je dovzetna za različne napade [6]. Veliki podatki (ang. big data) ponujajo epidemiologom, zdravnikom in strokovnjakom za zdravstveno politiko veliko priložnost za presojo na podlagi analize dostopnih podatkov, ki bo sčasoma dvignila raven oskrbe bolnikov [8].



Slika 1: Iskanje ravnovesja med zagotavljanjem zasebnosti uporabnikov in izvajanjem kvalitetnega nujenja podatkov o uporabnikih Vir: povzeto po [9]

Kot prikazuje slika 1, je izziv iskanje metodologij in tehnik za zaščito zasebnosti in nerazkrivanje občutljivih podatkov na eni strani in nudenje kvalitetnih podatkov o uporabnikih za raziskave, analizo in ustvarjanje dodane vrednosti z novimi znanji na drugi strani. K reševanju tega izziva je potrebno pristopiti skrbno, saj lahko poseganje v same podatke in njihovo preoblikovanje okrne njihovo uporabnost, kar lahko vodi v napačne interpretacije zavajajoče informacije ter neustrezne odločitve [2].

Za reševanje te vrste izzivov se je na področju podatkovnega rudarjenja razvilo posebno podpodročje, ki se ukvarja z ohranjanjem zasebnosti pri rudarjenju dokumentov z občutljivimi podatki. To področje se je v tuji literaturi uveljavilo pod imenom »Privacy preserving data mining« – PPDM [10], [11]. Metodologije PPDM so zasnovane tako, da zagotavljajo določeno raven zasebnosti, obenem pa ohranjajo uporabnost podatkov tako, da lahko rudarjenje še vedno učinkovito izvajamo na transformiranih podatkih. Da gre za zelo pomembno področje, je opozorila tudi Evropska komisija v »Mnenju številka 5« (EK, 2014). PPDM temelji na uporabi različnih metod, ki prispevajo k ohranjanju zasebnosti, kot npr. anonimizacija, randomizacija, uporaba permutacij, vnašanje šuma v podatke, kriptografske tehnike in druge [5]. Ker je področje razmeroma novo, a hkrati zanimivo, posebno danes, pri novih izzivih v zdravstvu (Covid19), bomo v nadaljevanju pripravili sistematičen pregled metod, ki se z izbrano problematiko ukvarjajo.

2. PODROČJA UPORABE PPDM

V tem razdelku bomo navedli področja, kjer prihaja do potreb po rudarjenju podatkov/dokumentov, ki

so po svoji naravi lahko občutljivi (Cranor et al., 2016). Čeprav avtorji navajajo domene uporab v ZDA, lahko izpostavljena področja uporabe preslikamo tudi v našo regijo. Področja, kjer je potreba po uporabi metod PPDM pogosta in intenzivna, so naslednja:

- **Zdravstveno varstvo.** Informatizacija procesov v zdravstvu lahko močno izboljša zdravstvene storitve, vključno z možnostjo natančnejše diagnostike, omogočanjem bolj prilagojene in usklajene oskrbe, hitrejšega razvoja novih načinov zdravljenja, učinkovitejšega zdravljenja ob nižjih stroških. Izziv predstavlja razkritje občutljivih zdravstvenih podatkov, širjenje le-teh (legalno ali nelegalno), kar lahko privede do različnih neprijetnih in uporabniku škodljivih situacij (npr. do diskriminacije pri zaposlovanju¹).
- **Transport.** Koristi informacijske tehnologije pri prevozu so lahko v zmanjšanju zastojev, preprečevanju nesreč, zmanjšanju smrti in poškodb, povečanju učinkovitosti porabe goriva, itd. Skrbi glede zasebnosti izvirajo iz možnosti sledenja gibanj posameznikov preko navigacijskih sistemov, cestnih senzorjev, prometnih kamer, zbiranja podatkov v avtomobilu in komunikacije med avtomobili.
- **Kazenski pregon in nacionalna varnost.** Organi pregona in obveščevalne agencije zbirajo in analizirajo različne vrste podatkov (kazenski zapisi in dopolnilne informacije) z namenom ustvarjanja »virtualne slike« posameznika, ki pomaga pri reševanju kriminalnih dejanj, preprečevanju napadov in sledenju teroristom. Glede ohranjanja zasebnosti so glavne skrbi te, da organizacije kot npr. policija množično zbirajo informacije o splošni populaciji, kar povečuje možnost nedovoljene uporabe in s tem niža kvaliteto učinka nadzora zaradi nezakonitega odtekanja podatkov. Kot primer lahko navedemo primer delovanja slovenske policije².
- **Izobraževanje.** Informacijska tehnologija in podatki o izobraževanju lahko izboljšajo izobraževa-

nje z nudenjem prilagodljivih in prilagojenih vsebin in spletnih tečajev. Tveganje glede zasebnosti izhaja iz občutljivosti podatkov o angažiranosti in uspešnosti uporabnikov (učenci, dijaki, študenti).

- **Sodobne internetne storitve.** Iskalniki, družbena omrežja, spletne video storitve in spletni trgovci imajo dostop do bogatega niza podatkov, ki jih je mogoče uporabiti za koristne namene, vključno z napredno personalizacijo vsebine in povezovanjem z drugimi (navadno poslovnimi) subjekti. Zaskrbljenost se kaže pri uporabi, zlorabi in skupni rabi podatkov za namene izven področja namembnosti.
- **Sodobni aplikacijski ekosistemi.** Naprave, kot so pametni telefoni, spletni brskalniki, pametne ure in njihove aplikacije, zagotavljajo uporabnikom veliko uporabnost (npr. pri športnih aktivnostih zaradi vgrajenih GPS naprav in pedimetrov), zabavo in funkcionalnost. Kot potencialno težavo lahko navedemo sledljivost uporabnika (zaradi GPS podatkov, ki se zbirajo v aplikaciji). Izziv predstavlja zagotovitev, da aplikacije spoštujejo zasebnost in varnost njihovih uporabnikov.
- **Internet stvari.** Pametna mesta, pametne zgradbe, pametni domovi, pametni hladilniki, televizije ipd. omogočajo izboljšanje življenjskih razmer, produktivnosti in kakovosti življenja. Vendar pa se lahko isti podatki uporabijo za sledenje, kdaj so posamezniki doma, katere TV programe gledajo, katera spletna mesta obiskujejo, njihovega urnika spanja in drugega vedenja. Tveganje predstavlja izkoriščanje takšnih podatkov za druge namene, kot npr. zavarovalne police (na voljo so informacije o prehrabnih navadah, aktivnostih in s tem tveganjih – profiliranje uporabnikovih navad), nezaželeno oglaševanje ali kriminal.
- **Finančni sektor.** Podatki finančnih institucij lahko regulatorjem pomagajo pri oceni skladnosti in omogočijo analizo trendov ter opozorijo na nevarnosti kot npr. prihajajoča finančna kriza. Vendar so finančni podatki občutljivi ne le na ravni posameznih strank, temveč tudi na ravni institucij, saj razkrivajo lastniške informacije o strategijah in tržnih deležih.
- **Odpri podatki državne uprave.** Vlade na vseh ravneh sproščajo velike količine podatkov, da bi povečale zaupanje in preglednost ter omogočile inovativne aplikacije. Vendar se te objave podatkov pogosto nanašajo na občutljive informacije o

¹ Tri tedne po tem, ko je Nydia Velázquez zmagala, kot kandidatka Demokratske stranke v New Yorku za predstavnico v Ameriškem predstavnem domu, je nekdo iz bolnišnice St. Claire v New Yorku, preko faksa poslal Velázquezino zdravstveno kartoteko časopisu New York Post. V dokumentu je bila podrobno opisana oskrba pacientke, ki je tam pristala zaradi poskusa samomora. Poskus samomora se je zgodil nekaj let pred volitvami, na katerih je zmagala. Povzeto po (Wu and Velázquez, 2000).

² Avstrijska nevladna organizacija AlgorithmWatch je v svojem poročilu analize policijske uporabe tehnologije za prepoznavanje obrazov zapisala, da slovenska policija od leta 2014 uporablja doma razvito tehnologijo za prepoznavanje obrazov. Kot navajajo, gre za problem regulacije te tehnologije. Informacijska pooblaščenka je tako med leti 2015 in 2019 izdala več negativnih mnenj Ministrstvu za notranje zadeve. Slovenska policija in biometrijske metode nadzora <https://www.eticen.it/2019/12/12/slovenska-policija-in-biometrijske-metode-nadzora>, dostopano: januar 2020

državljanih. Lahko navedemo nekaj primerov takšnih spletnih storitev pri nas: eDavki³, eUprava⁴, eVem⁵, eZdravje⁶ in drugi.

3. METODE ZA ZAŠČITO OBČUTLJIVIH PODATKOV

Analitika velikih podatkov (ang. big data) sestoji iz petih stopenj oz. faz, in sicer: pridobivanje podatkov, shranjevanje podatkov, upravljanje s podatki, analiza podatkov ter vizualizacija podatkov in poročanje. Pri dveh od teh se soočamo z ohranjanjem zasebnosti: **shranjevanje podatkov** in **upravljanje s podatki** [8], [13], [14]. Podatki, ki jih pridobivamo, so lahko strukturirani, delno strukturirani, ali pa gre za nestrukturirane podatke. Podatke lahko pridobimo iz ustnih virov (intervju, telefonski pogovor) ali pisnih virov (npr. anketa, izvid, vprašalnik, diagnoza, mnenje). Pogosto so viri podatkov tudi slikovni ali multimedijški (npr. magnetna resonanca, računalniška tomografija, ultrazvok itd.). Ne nazadnje, v dobi interneta lahko podatke pridobimo tudi iz spletnih anket, vprašalnikov, poskusov. Pogosto podatke že sami shranjujemo, posredujemo in s tem ponujamo v varne ali ne – oblačne storitve (ang. cloud services). Gre za podatke, pridobljene iz pametnih naprav, telefonov z uporabo aplikacij, kot so npr. Drive, Training Peaks, Polar Flow, Strava, in podobne.

Glavni nalogi uporabe PDDM, kot ju navedejo Xu, Jiang, Wang, Yuan, & Ren (2014), sta soočanje in razreševanje problematike neprimernosti neposredne uporabe občutljivih, surovih podatkov (npr. številka osebne izkaznice, mobilnega telefona) za rudarjenje ter potreba po izključitvi občutljivih rezultatov rudarjenja, katerih razkritje bi povzročilo kršitev zasebnosti. Pionirsko delo, opis prvih metod rudarjenja občutljivih podatkov na tem področju najdemo v člankih [16], [17].

Pri rudarjenju podatkov z namenom zaščite zasebnosti se uporabljajo različne metode. Usmerjene so predvsem v omejevanje dostopa in uporabe občutljivih podatkov, ki bi sicer lahko identificirali posameznika, za nadaljnjo analizo. Po Abdul, Aldeen, Salleh, & Razzaque (2015), Qi & Zong (2012) in Taneja, Khanna, & Tilwalia, (2016) med glavne metode umeščamo:

- **Razprševanje podatkov** (ang. partitioning): podatke distribuiramo po eni ali več podatkovnih baz.
- **Izkrivljanje podatkov** (ang. data distortion, perturbation): pri tem načinu posegamo v podatke, ki jih želimo uporabiti in katerih vrednost je tista, ki jo želimo zaščititi. Sem umeščamo **anonimizacijo**, **randomizacijo**, **vrtenje**, **vnašanje šuma** v podatke.
- **Kriptografske tehnike šifriranja** (ang. cryptographic technique): gre za različne, pogosto računsko potratne metode, kjer se podatki s pomočjo ustreznega šifrirnega algoritma (simetričnega ali asimetričnega) in šifrirnega ključa pretvorijo v neberljivo obliko.

V nadaljevanju bomo posamezne metode nekoliko podrobneje opisali.

3.1 Razprševanje podatkov

O razprševanju podatkov (ang. partitioning) govorimo, kadar podatke, shranjene v eni podatkovni bazi, porazdelimo (razpršimo) v več podatkovnih baz. Podatke lahko razpršimo horizontalno, vertikalno ali funkcionalno. Vse omenjeno lahko počnemo centralizirano (na enem mestu) ali pa distribuirano (na več lokacijah).

Pri horizontalni razpršitvi podatkov pridobimo razširljivost (ang. scalability) in učinkovitosti v smislu hitrejšega dostopa do podatkov (ang. performance), na varnosti (ang. security) pa precej manj, saj so v posamezni relaciji (povezava med dvema ali več entitetami) zbrani vsi atributi. Iz vidika varnosti je poskrbljeno zgolj za to, da niso vsi podatki o vseh subjektih zbrani na enem mestu, ampak so razpršeni po več podatkovnih bazah.

Pri vertikalni razpršitvi gre za nasproten proces. Hitrost dostopa in uporabe podatkov pada, saj je treba attribute iz več podatkovnih baz med sabo združiti. Pri tem pristopu imamo attribute razdeljene v več skupin, vsaka izmed skupin pa je v svoji podatkovni bazi. Navadno ima vsaka relacija skupni ključ, ki povezuje podatke med seboj.

Pri funkcionalni razpršitvi ločujemo podatke glede na funkcijo, oz. uporabo⁷ [18], [21].

Omeniti je treba, da z razprševanjem samih podatkov ne spreminjamo, pač pa lahko zgolj ome-

³ Edavki, <https://edavki.durs.si/EdavkiPortal/OpenPortal/CommonPages/Opdypn/PageA.aspx>

⁴ eUprava, <https://e-uprava.gov.si>

⁵ eVem, <http://evem.gov.si/evem/drzavljanizacetna.evem>

⁶ eZdravje, <https://zvem.ezdrav.si/e-zdravje>

⁷ Horizontal, vertical, and functional data partitioning, <https://docs.microsoft.com/en-us/azure/architecture/best-practices/data-partitioning>, dostopano januar 2010

jimo dostop do njih. Pri vertikalni razpršitvi lahko uporabniku ponudimo dostop do delov podatkov (atributov, ki jih potrebuje), uporabnik pa si ne more ustvariti »celotne slike«, ker je dostop do celotne vsebine omejen. Za prikaz celotne slike moramo združiti podatke iz različnih podatkovnih baz oz. sistemov.

3.2 Izkrivljanje podatkov

V skupino izkrivljanja podatkov (ang. data perturbation) sodijo metode, tehnike in algoritmi, ki podatke spreminjajo, ali pa vanje dodajajo šum. Največ literature s področja PPDM je posvečeno ravno temu segmentu, ki je tudi najbolj kompleksen. Večinoma gre za matematične metode, ki posegajo v podatke in jih z uporabo vektorjev, matrik, faktorjev – spreminjajo. Na začetku je bilo implementiranih nekaj metod, ki so temeljile zgolj na naključnem seštevanju in množenju, a le-te niso bile imune na praktično nobeno vrsto napada [22].

Cilj izkrivljanja podatkov je ponuditi informacije, ki jih je mogoče uporabiti za rudarjenje na način, da ostane prikrita identifikacija lastnika (subjekta) atributov. Attribute v grobem delimo v tri skupine: *identifikacijski atributi* (identificirajo subjekt), *javni* ali *kvazi atributi* (njihove vrednosti lahko pridobimo tudi v drugih, javnih bazah podatkov, kot je npr. volilni imenik, podatki v profilu na socialnih omrežjih (letnik ali starost, kraj, naslov, itd.), ter *privatni atributi*, ki npr. v primeru bolnišničnega kartona opisujejo stanje bolnika, bolezen, zdravljenje. Cilj je zagotoviti dostop do privatnih atributov, s pomočjo katerih bi radi ugotovili povezave, odvisnosti in s tem prišli do novih spoznanj in znanja, hkrati pa zavarovali informacije, ki identificirajo posameznika [23].

Medtem, ko pri metodah razprševanja in šifriranja podatke skrivamo, delimo, distribuiramo, jih s pomočjo anonimizacijskih tehnik **izkrivljamo** – z namenom nudenja nadaljnji uporabi. Obstoječe metode na področju izkrivljanja podatkov opisujejo Sachan, Roy, & Arun (2013) in Qi & Zong (2012), ki med metodami omenjajo metodo **k-anonimnosti**, **generalizacijo**, **klasifikacijo**, **gručenje**, **povezovalna pravila**, **porazdeljeno ohranjanje zasebnosti**, **l-raznolikost**, **t-podobnost**, **randomizacijo** in **drevesno razvrščanje**.

Metode, ki uporabljajo algoritme dodajanja šuma, permutacij in randomizacijske tehnike, imajo prednosti, kot so npr. neodvisno izvajanje skozi vse vrednosti atributov (neodvisnost) in ohranjanje statistične natančnosti po rekonstrukciji originalnih podat-

kov. Med slabosti pa umeščamo zmanjšano uporabnost atributov pri generalizaciji v intervale, skrivanje robnih podatkov, kar zahteva visoko uporabo šuma in s tem znižuje uporabnost podatkov [2]. V članku Li & Sarkar (2006) ponujata izboljššan pristop spreminjanja podatkov (dodajanja šuma), kot je zgolj povečevanje / zmanjševanje numeričnih vrednosti atributa za isti faktor, vrednost ali rotacijo. Ta temelji na razvrščanju v drevesa in na uporabi največje variance vrednosti med atributi.

Generalizacijo uporabljamo v primeru, ko uporabnik podatkov ne potrebuje natančnih vrednosti atributov in lahko le-te posplošimo, npr. v intervale, ki niso nujno vedno enako široki. Pri osebni dohodku na primer lahko za nižje vrednosti uporabimo ožji, pri višjih vrednostih pa širši interval. Namesto prave vrednosti ponudimo interval, na katerem vrednost leži. Temu načinu spreminjanja rečemo *interval vrednostnih razredov* (ang. value-class interval). Drugi način je *vrednostno izkrivljanje* (ang. value distortion). Gre za dodajanje šuma, saj namesto vrednosti x_i ponudimo spremenjeno vrednost $Z_{(i)} = x_i + r$, kjer je r naključna vrednost iz nekega intervala $[-a, +a]$, ali pa generirana s pomočjo normalne (Gaussove) porazdelitve [10]. Pri izkrivljanju podatkov s to metodo, govorimo o izkrivljanju podatkov z dodatnim šumom, multiplikativnim šumom ali permutacijo [22], [26]. Kot šum si lahko predstavljamo podatek, ki je spremenjen na nivoju vseh elementov (npr. podatek o teži, višini pacientov, je za vse paciente spremenjen za določen faktor).

k-anonimizacijo uvrščamo v drugo skupino anonimizacijskih tehnik. Cilj te tehnike je anonimizirati člana množice ali skupine na način generalizacije vrednosti atributov (npr. mesto ali poštno številko z regijo, višino zaokrožiti na desetice, starost v interval, itd.). Pristop k-anonimnosti sta prva predlagala Samarati & Sweeney v (1998) in Sweeney (2002). Cilj postopka k-anonimizacije je vključitev vsakega posameznika, o katerem podatki so nam na voljo, v večjo skupino, s k-posamezniki in s tem povečati negotovost identifikacije posameznika. Tehnika k-anonimizacija ni imuna na uporabo posrednega znanja (ang. background knowledge), ki ga ima napadalec. To nastane zaradi slabo definiranih intervalov ali znanja, ki ga lahko ima napadalec o posamezniku, katerega podatke preiskuje (npr. zaradi atributov, ki niso skriti). Mendes & Vilela (2017) navajata, da sta prednosti metode k-anonimizacija enostavnost defi-

Tabela 1: Originalna tabela s podatki o pacientih Vir: Povzeto po Lin (2016)

Identif. atributi		Atributi kvazi identifikatorji (QI)			Občutljivi atributi
ID	Ime	Starost	Holesterol	Trigliceridi	Bolezen
45435	Janez	22	6	1.8	bolezen srca
46434	Petra	26	5.9	1.7	rak
65675	Tilen	36	7	2.3	hipertenzija
34567	Meta	35	8	3.9	rak
54345	Andrej	49	6.4	4.3	bolezen srca
34333	Špela	44	5.5	5.9	hipertenzija

nicije protokola in velik nabor obstoječih algoritmov za doseg k-anonimizacije, kot slabost pa predvidevanje, da vsak zapis predstavlja podatke o edinstvenem posamezniku. Če ni tako, se razred enakovrednosti s k zapisi ne poveže nujno s k različnimi posamezniki. Prav tako privatni (občutljivi) atributi ne pridejo v poštev za anonimizacijo v primeru, če imajo vsi podatki razreda isto vrednost.

V tabeli 1 in 2 prikazujemo primer originalne in spremenjene tabele s podatki o pacientih.

V tabeli 1 vidimo namišljene podatke o bolnikih, razdeljene v tri skupine atributov (identifikacijski, kvazi identifikatorji ter občutljivi ali privatni). Podatke bi radi anonimizirali do te mere, da ne bi bilo možno neposredno povezati bolnika z boleznijo, hkrati pa bi zaradi nadaljnje analize radi obdržali čim več podatkov, zanimivih za raziskave in razvoj znanosti na področju medicine. Rezultate modifikacije podatkov prikazujeta tabeli 2 in 3.

V tabeli 2 smo odstranili identifikacijske attribute in generalizirali tri druge. Prikazan generalizacijski postopek imenujemo k-anonimizacija [23], [26], [29], [30]. Cilj postopka k-anonimizacije je vključitev vsakega posameznika, o katerem podatki so nam na vo-

ljo, v večjo skupino s k -posamezniki in s tem povečati negotovost identifikacije posameznika. V opisanem primeru (v tabelah 1 in 2) prikazujemo 2-anonimnost, saj za vsako kombinacijo generaliziranih kvazi atributov, obstajata vsaj dve bolezni.

Ena od slabosti opisanega postopka je premajhna vrednost k in dejstvo, da tehnika ne uporablja dodajanja šuma, kot npr. randomizacija. S sklepanjem ali ugibanjem je pri majhnem k verjetnost, da imata osebi eno ali drugo bolezen večja, kot pri velikem k . Npr., če vemo, da so v razpredelnici podatki nekoga, za katelega približno poznamo kvazi identifikatorje, vemo, da ima bodisi eno ali pa drugo bolezen (pri 2-anonimnosti). V članku Mendes & Vilela (2017), avtorja med prednosti metode k-anonimizacija umeščata:

- enostavnost definicije protokola,
- velik nabor obstoječih algoritmov za doseg k-anonimizacije, med slabosti pa:
- Predvideva se, da vsak zapis predstavlja podatke o edinstvenem posamezniku. Če ni tako, se razred enakovrednosti s k zapisi ne poveže nujno s k različnimi posamezniki.
- Privatni (občutljivi) atributi ne pridejo v poštev za anonimizacijo, kar lahko prispeva k razkritju in-

Tabela 2: Spremenjena tabela s podatki o pacientih Vir: Povzeto po Lin (2016)

Atributi kvazi identifikatorjev (QI)			Občutljivi atributi
Starost	Holesterol	Trigliceridi	Bolezen
[20-29]	[5-7]	[0-2]	bolezen srca
[20-29]	[5-7]	[0-2]	rak
[30-39]	[7-9]	[2-4]	hipertenzija
[30-39]	[7-9]	[2-4]	rak
[40-49]	[5-7]	[4-6]	bolezen srca
[40-49]	[5-7]	[4-6]	hipertenzija

formacije/podatka v primeru, če imajo vsi podatki razreda isto vrednost.

Uporaba je razširjena predvsem pri podatkih, ustvarjenih v zdravstvu in podatkih, ki vključujejo (geo)lokacijske informacije.

Nadgradnja k-anonimnosti je **l-raznolikost**. Dodana je še ena omejitev, in sicer, da se vsak atribut v ekvivalenčnem razredu pojavi vsaj l-krat tako, da je napadalec vedno precej negotov glede atributov tudi, če ima osnovne informacije o določenem posamezniku, na katerega se nanašajo osebni podatki (Machanavajhala, Gehrke, Kifer, & Venkatasubramaniam, 2006). Med slabostmi metode je treba poudariti predvsem dejstvo, da je zahtevna za implementacijo (težko doseči primerno obliko), poleg tega pa se napadalec, v primeru, da so občutljivi atributi nekega razreda enaki, nauči / izve vrednost tega atributa za določenega posameznika [2], [32].

V tabeli 3 prikazujemo primer l-raznolikosti, kjer je za ceno varnosti zmanjšana zrnatost podatkov. Podatki so generalizirani do te mere, da je ustvarjena 3-raznolikost, ker obstajajo trije različni občutljivi atributi znotraj vsakega bloka (bloka sta dva) v tabeli. S to potezo se zmanjša tveganje identifikacije, a na drugi strani bolj generalizira vrednosti kvazi identifikatorjev (povzeto po Machanavajhala, Gehrke, Kifer, & Venkatasubramaniam (2006).

Čeprav model l-raznolikost učinkovito rešuje težave, ki obstajajo v modelu k-anonimnosti, se model ne more »upreti« napadom na podobnost (ang. similarity attacks). To pomeni, da je delež vrednosti občutljivega atributa prevelik. V tem primeru je velika verjetnost, da bo napadalec razkril zasebnost posameznika. Zato je znanstvenik Li Ning Hui predlagal model t-podobnost (ang. t-closeness). Ta zahteva, da

vrednost razlike med porazdelitvijo vrednosti občutljivih atributov v enakovrednih razredih in porazdelitvijo atributa v celotni podatkovni tabeli ni večja od t. Če je na primer občutljivi atribut številski, l-raznolikost ne upošteva, da so si nekatere vrednosti lahko zelo podobne (da so si blizu), kar rešuje metoda t-podobnost. Ta določa, da mora biti porazdelitev občutljivega atributa v vsakem ekvivalenčnem razredu podobna porazdelitvi v celotni tabeli. To lahko prepreči napade na podobnost in dodatno reši težave, ki obstajajo v modelu l-raznolikosti. Model t-podobnost je velja za najboljši anonimizacijski model varovanja zasebnosti [26], [31], [33], [34].

3.3 Šifriranje podatkov

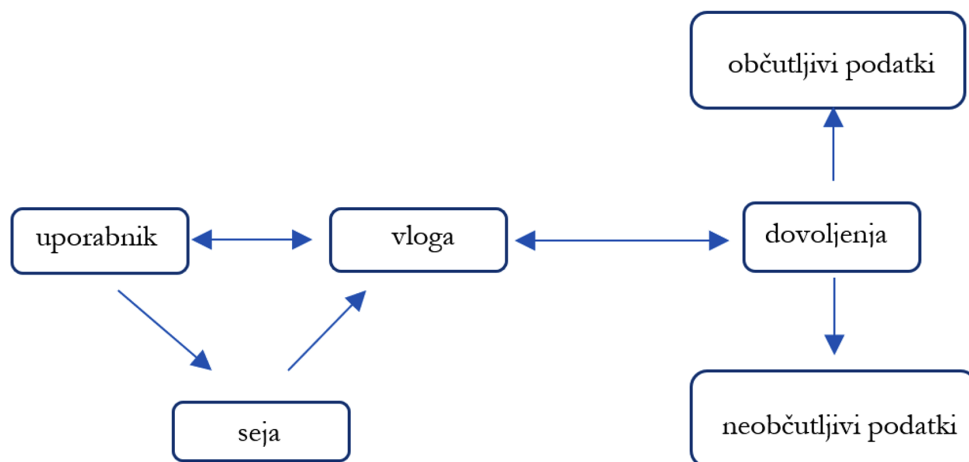
Šifriranje občutljivih podatkov (atributov) je dober pristop k procesu varovanju podatkov, saj ne spremeni podatkov, ne prihaja do izgube (generalizacije), ali šuma. Med slabosti štejemo predvsem težavno implementacijo pri velikih zbirkah podatkov, poleg tega pa rezultat (originalni podatki) odkriva tako javne kot skrite attribute [20], [35].

Metode šifriranja podatkov lahko uporabimo v kombinaciji z razprševanjem podatkov, in sicer na dva načina: podatki so lahko razdeljeni vertikalno skozi več podatkovnih baz (med več lastniki), kjer so šifrirani zgolj zasebni podatki, lahko pa so šifrirani vsi podatki. V praksi se je razvil model »ohranjanje varnosti na podlagi kontrole dostopa preko vlog« (ang. privacy preserving role based access control approach – PRBAC), ki je ponazorjen na sliki 5. Ta model kombinira vertikalno razprševanje podatkov in tehnologijo šifriranja za dostop do podatkov, ki jih deli na javne in zasebne [36].

PRBAC je eden izmed pristopov za zaščito informacij v relacijski bazi podatkov, ki uporabnikom

Tabela 3: dodatno spremenjena tabela s podatki o pacientih Vir: Povzeto po Lin (2016)

Atributi kvazi identifikatorjev (QI)			Občutljivi atributi
Starost	Holesterol	Trigliceridi	Bolezen
< 50	[5-7]	[0-6]	bolezen srca
< 50	[5-7]	[0-6]	rak
< 50	[5-7]	[0-6]	hipertenzija
< 50	[5-9]	[2-6]	hipertenzija
< 50	[5-9]	[2-6]	rak
< 50	[5-9]	[2-6]	bolezen srca



Slika 5: Model PRBAC Vir: [36]

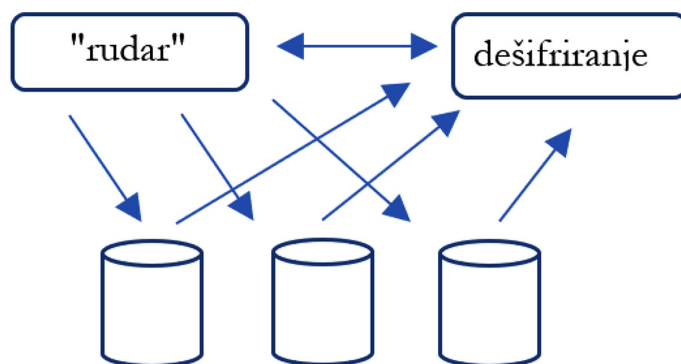
preprečuje pridobivanje dovolj velikih in raznolikih vzorcev baze podatkov. Prav tako onemogoča sledenje vzorcem, ki jih ne bi smeli razkriti. To omogoča razvrščanje podatkov na javne (neobčutljive) in tajne (občutljive). Slabost metode je časovna potratnost (zaradi zahtevnosti postopkov šifriranja in dešifriranja) ter čas, potreben za sestavljanje relacije (iz n podatkovnih baz). Postopek je moč pohitriti z uvedbo horizontalnega razprševanja podatkov. Tedaj je iskana relacija v celoti v eni izmed podatkovnih baz, a je takšna realizacija manj varna [36].

Pri uporabi *razprševanja podatkov* in šifriranja, ali uporabi kombinacije teh dveh metod, podatki ostajajo v obliki, kot so nastali. S tem, ko jih ne spreminjamo ali izkrivljamo, ostajajo za nadaljnjo analizo potencialno najustrežnejši, a bi pri uporabi lahko prišlo do kršitve ohranjanja zasebnosti, predvsem v primeru, če bi do podatkov lahko dostopali nelegalno ali nepooblaščno.

Model PRBAC in vertikalno, ter funkcionalno porazdeljene podatkovne baze težavo delno odpravljajo, saj je v prvem primeru zagotovljen dostop do podatkov preko vlog dostopa, v drugem pa do pridobitve zgolj dela podatkov (dela atributov relacije). Nobeden od teh dveh pristopov ni primeren za nudenje podatkov v množično uporabo (npr. za podatkovno rudarjenje), saj ni poskrbljeno za anonimizacijo podatkov oz. za zakritje povezave med lastnikom podatkom in njegovimi atributi v relaciji drugače, kot s ponujanjem dostopa do zgolj dela podatkov končnemu uporabniku (slika 6).

4. ZAKLJUČEK

Pri iskanju ravnovesja med nudenjem podatkov za analizo, z namenom ustvarjanja dodane vrednosti in znanj, je potrebno poskrbeti za zaščito identitete in interesov posameznika po anonimnosti. Pri pregledu literature smo zasledili, da se za reševanje omenjenih



Slika 5: Dešifriranje vertikalno porazdeljenih podatkov Vir: [36]

izzivov na področju podatkovnega rudarjenja razvilo posebno podpodročje, ki se ukvarja z ohranjanjem zasebnosti pri tem procesu. Metode, s katerimi dosegamo zaščito posameznika na eni in uporabno vrednost podatkov na drugi strani v grobem delimo na metode *razprševanja podatkov*, *metode izkrivljanja in metode šifriranja podatkov*. Za doseganje višje zaščite lahko uporabimo tudi kombinacije teh metod.

Za primere iz prakse, kjer bi uporaba tovrstnih metod pripomogla k boljši uporabi (javnih) podatkov in informacij javnega značaja, omenimo Zakon o dostopu do informacij javnega značaja (ZDIJZ-NPB10)⁸. Ta v šestem členu navaja izjeme, kjer dostop do takšnih podatkov ni dovoljen, večinoma zaradi posledic razkritja in s tem kršitve varstva osebnih, ali drugih podatkov. Pri Splošni uredbi o varstvu podatkov (ang. general data protection regulation – GDPR⁹), uredba v 17. členu določa Pravico do izbrisa (»pravico do pozabe«), ki določa pogoje, pod katerimi lahko posameznik zahteva izbris osebnih podatkov iz dokumentov, ter 28. člen, kjer določa pristojnosti in omejitve obdelovalca podatkov. Za omogočanje dostopa do dokumentov takšne narave sta predlagana šifriranje in revizija dostopov [37] ter izkrivljanje podatkov. Skladno z GDPR se lahko dokumenti objavijo po odstranitvi vseh podatkov, ki identificirajo posameznika [38].

Največ izboljšav in novih metod in algoritmov smo zasledili ravno za področje spreminjanja podatkov (ang. data perturbation). To je edini pristop, ki podatke transformira, a pri tem ohranja (skozi funkcijo transformacije) možnost rekonstrukcije, podatki pa ohranjajo visoko uporabno vrednost za nadaljnjo uporabo – podatkovno rudarjenje.

VIRI IN LITERATURA

- [1] J. R. Smith and S.-F. Chang, »New visual information in the form of images Visually Searching the Web for Content,« pp. 12–20, 1997, doi: 10.1080/10413200.2012.704621.
- [2] R. Mendes and J. P. Vilela, »Privacy-Preserving Data Mining: Methods, Metrics, and Applications,« *IEEE Access*, vol. 5, pp. 10562–10582, 2017, doi: 10.1109/ACCESS.2017.2706947.
- [3] R. Devakunchari, »Analysis on big data over the years,« *International Journal of Scientific and Research Publications*, vol. 4, no. 1, pp. 1–7, 2014, [Online]. Available: www.ijsrp.org.
- [4] H. Sahu, S. Shirma, and S. Gondhalakar, »A Brief Overview on Data Mining Survey,« *Ijctee*, vol. 1, no. 3, pp. 114–121, 2008.
- [5] A. S. Shanthi and M. Karthikeyan, »A review on privacy preserving data mining,« *2012 IEEE International Conference on Computational Intelligence and Computing Research, ICCIC 2012*, 2012, doi: 10.1109/ICCIC.2012.6510302.
- [6] A. Singh, »Data Publishing Techniques and Privacy Preserving,« *Int. J. Inf. Secur. Res.*, vol. 9, no. 3, pp. 1–23, 2019.
- [7] K. (Eds. . Markov, A., Polak Petrič, A., & Nastovska, *Splošna deklaracija človekovih pravic*. Ljubljana : Fakulteta za družbene vede, Založba FDV : Ministrstvo za zunanje zadeve Republike Slovenije, 2018, 2018.
- [8] S. SA, »Big Data in Healthcare Management: A Review of Literature,« *American Journal of Theoretical and Applied Business*, vol. 4, no. 2, p. 57, 2018, doi: 10.11648/j.ajtab.20180402.14.
- [9] L. Cranor, T. Rabin, V. Shmatikov, S. Vadhan, and D. Weitzner, »Towards a Privacy Research Roadmap for the Computing Community,« *Http://Cra.Org/Ccc/Resources/Ccc-Led-Whitepapers/*. pp. 1–23, 2016, [Online]. Available: http://arxiv.org/abs/1604.03160.
- [10] R. Agrawal and R. Srikant, »Privacy-preserving data mining,« *SIGMOD Rec. (ACM Spec. Interes. Gr. Manag. Data)*, vol. 29, no. 2, pp. 439–450, Jan. 2000, doi: 10.1145/335191.335438.
- [11] Y. Lindell and B. Pinkas, »Privacy Preserving Data Mining,« pp. 36–54, 2000.
- [12] L. Cranor, T. Rabin, V. Shmatikov, S. Vadhan, and D. Weitzner, »Towards a Privacy Research Roadmap for the Computing Community,« *Http://Cra.Org/Ccc/Resources/Ccc-Led-Whitepapers/*. pp. 1–23, 2016.
- [13] K. Vassakis, E. Petrakis, and I. Kopanakis, »Big Data Analytics: Applications, Prospects and Challenges,« 2018, pp. 3–20.
- [14] A. Pawar, S. Ahirrao, and P. P. Churi, »Anonymization Techniques for Protecting Privacy : A Survey,«
- [15] L. Xu, C. Jiang, J. Wang, J. Yuan, and Y. Ren, »Information security in big data: Privacy and data mining,« *IEEE Access*, vol. 2, no. January, pp. 1151–1178, 2014, doi: 10.1109/ACCESS.2014.2362522.
- [16] R. Agrawal and R. Srikant, »Privacy-preserving data mining,« *SIGMOD Rec. (ACM Spec. Interes. Gr. Manag. Data)*, vol. 29, no. 2, pp. 439–450, Jan. 2000, doi: 10.1145/335191.335438.
- [17] Y. Lindell and B. Pinkas, »Privacy Preserving Data Mining BT – Advances in Cryptology — CRYPTO 2000,« 2000, pp. 36–54.
- [18] X. Qi and M. Zong, »An Overview of Privacy Preserving Data Mining,« *Procedia Environ. Sci.*, vol. 12, no. Part B, pp. 1341–1347, Jan. 2012, [Online]. Available: http://10.0.3.248/j.proenv.2012.01.432.
- [19] Y. Abdul, A. S. Aldeen, M. Salleh, and M. A. Razzaque, »A comprehensive review on privacy preserving data mining,« *Springerplus*, 2015, doi: 10.1186/s40064-015-1481-x.
- [20] S. Taneja, S. Khanna, and S. Tilwalia, »A Review on Privacy Preserving Data Mining: Techniques and Research Challenges,« vol. 6, no. 3, pp. 35–40, 2016.
- [21] P. P. Panse and P. L. Paikrao, »Survey of Privacy Preserving Techniques and Upcoming Techniques : A Review,« vol. 6, no. 2, pp. 1798–1802, 2017.
- [22] S. Upadhyay, C. Sharma, P. Sharma, P. Bharadwaj, and K. R. Seeja, »Privacy preserving data mining with 3-D rotation transformation,« *J. King Saud Univ. – Comput. Inf. Sci.*, vol. 30, no. 4, pp. 524–530, 2018, doi: 10.1016/j.jksuci.2016.11.009.
- [23] A. Gionis and T. Tassa, »K-anonymization with minimal loss of information,« *IEEE Transactions on Knowledge and Data Engineering*, vol. 21, no. 2, pp. 206–219, 2009, doi: 10.1109/TKDE.2008.129.
- [24] A. Sachan, D. Roy, and P. V Arun, »An Analysis of Privacy Preservation Techniques in Data Mining BT – Advances in Computing and Information Technology,« 2013, pp. 119–128.

⁸ Zakon o dostopu do informacij javnega značaja, <https://zakonodaja.com/zakon/zdijz>, dostopano september 2021

⁹ GDPR, <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:02016R0679-20160504>, dostopano september 2021

- [25] X. B. Li and S. Sarkar, »A tree-based data perturbation approach for privacy-preserving data mining,« *IEEE Transactions on Knowledge and Data Engineering*, vol. 18, no. 9. pp. 1278–1283, 2006, doi: 10.1109/TKDE.2006.136.
- [26] EK, »Mnenje št. 5/2014 o anonimizacijskih tehnikah.« Evropska komisija, Delovna skupina za varstvo podatkov člana 29, 2014, [Online]. Available: https://ec.europa.eu/justice/article-29/documentation/opinion-recommendation/files/2014/wp216_sl.pdf.
- [27] P. Samarati and L. Sweeney, »Generalizing data to provide anonymity when disclosing information (abstract),« 1998.
- [28] L. Sweeney, »K-Anonymity: A Model for Protecting Privacy,« *Int. J. Uncertain. Fuzziness Knowl.-Based Syst.*, vol. 10, no. 5, pp. 557–570, 2002, doi: 10.1142/S0218488502001648.
- [29] C. Y. Lin, »A reversible data transform algorithm using integer transform for privacy-preserving data mining,« *J. Syst. Softw.*, vol. 117, pp. 104–112, 2016, doi: 10.1016/j.jss.2016.02.005.
- [30] A. Gkoulalas-Divanis and V. S. Verykiosc, »An overview of privacy preserving data mining,« *Crossroads*, vol. 15, no. 4, pp. 23–26, 2009, doi: 10.1145/1558897.1558903.
- [31] A. Machanavajjhala, J. Gehrke, D. Kifer, and M. Venkatasubramaniam, » l -Diversity: Privacy beyond k -anonymity,« *Proceedings – International Conference on Data Engineering*, vol. 2006. p. 24, 2006, doi: 10.1109/ICDE.2006.1.
- [32] J. Vasa and P. Modi, »Review of Different Privacy Preserving Techniques in PPDP,« *International Journal of Engineering Trends and Technology*, vol. 59, no. 5. pp. 223–227, 2018, doi: 10.14445/22315381/ijett-v59p242.
- [33] G. Hao and X. Ya-Bin, »Research on privacy preserving method based on T-closeness model,« in *2017 3rd IEEE International Conference on Computer and Communications (ICCC)*, 2017, pp. 1455–1459, doi: 10.1109/CompComm.2017.8322783.
- [34] P. (1) Quirós, P. (1) Alonso, I. (2) Díaz, and S. (3) Montes, »Protecting data: a fuzzy approach,« *Int. J. Comput. Math.*, vol. 92, no. 9, pp. 1989–2000, Sep. 2015, doi: 10.1080/00207160.2014.928700.
- [35] B. Pinkas, »Cryptographic Techniques for Privacy-Preserving Data Mining,« *SIGKDD Explor. Newsl.*, vol. 4, no. 2, pp. 12–19, Dec. 2002, doi: 10.1145/772862.772865.
- [36] L. Vasudevan, S. E. D. Sukanya, and N. Aarathi, »Privacy preserving data mining using cryptographic role based access control approach,« *Imecs 2008: International Multiconference of Engineers and Computer Scientists, Vols I and II*. pp. 474–479, 2008.
- [37] C. Evans, »HOW GDPR WILL SHAKE UP DATA STORAGE.,« *Comput. Wkly.*, pp. 25–28, Aug. 2017, [Online]. Available: <http://nukweb.nuk.uni-lj.si/login?url=https://search.ebscohost.com/login.aspx?direct=true&AuthType=ip,uid&db=f5h&AN=124782283&lang=sl&site=eds-live&scope=site>.
- [38] K. Broen, R. Trangucci, and J. Zelner, »Measuring the impact of spatial perturbations on the relationship between data privacy and validity of descriptive statistics,« *Int. J. Health Geogr.*, vol. 20, no. 1, pp. 1–17, 2021, doi: 10.1186/s12942-020-00256-8.

■

Matjaž Kragelj Zaposlen v Narodni in univerzitetni knjižnici, zadnje desetletje kot vodja enote za informacijsko tehnologijo in digitalno knjižnico. Njegovo področje dela je stičišče računalništva in knjižnične stroke. Večina njegovih dejavnosti vključuje načrtovanje in upravljanje digitalne knjižnice in spletnih storitev, razvoj in usklajevanje dejavnosti nacionalnega agregatorja e-vsebin na področju kulture, izobraževanja in svetovanja na področju integriranega upravljanja digitalnih virov, sodelovanje pri razvoju in vzdrževanju celovitih informacijskih rešitev v knjižnici. Sodeluje tudi pri razvoju in vzdrževanju digitalnega arhiva knjižnice in storitev zajema in ohranjanja slovenskega spleta. Sodeluje pri digitalizaciji nacionalne kulturne dediščine, shranjene v knjižnici, ter na področju spremljanja in oblikovanja (mednarodnih) priporočil in smernic na področju pridobivanja, dolgoročnega arhiviranja in dostopa do digitalnih virov. Sodeloval je v več mednarodnih projektih in bil avtor več člankov s svojega področja dela. V zadnjih letih je več pozornosti namenil vizualizaciji podatkov in rudarjenja besedil, iskanju povezav in vzorcev v podatkih.

■

Mirjana Kljajić Borštnar je izredna profesorica za področje informacijskih sistemov na Fakulteti za organizacijske vede, Univerze v Mariboru. Njeno raziskovalno delo je usmerjeno v sisteme za podporo odločanju, odkrivanje znanja v podatkih in organizacijsko učenje. Izsledke raziskav objavlja v mednarodnih znanstvenih revijah in konferencah, med drugim *Expert Systems with Application*, *PLOS ONE*, *Industrial Management & Data Systems*, *System Dynamics Review*. Sodeluje v evropskih in domačih projektih. Je sovodja programskega odbora Blejske e-konference in Simpozija o operacijskih raziskavah v Sloveniji ter članica programskih odborov konferenc *DSI*, *DataScience*, *WorldCist* in drugih. V domačem okolju je aktivna kot predstavnica raziskovalnih organizacij v *SRIP PMIS* za področje *Ai*, *HPC & Big Data*, članica izvršnega odbora pobude *AI4Slovenia* in članica uredniškega odbora revije *Uporabna informatika*.

■

Alenka Brezavšček je izredni profesor na Fakulteti za organizacijske vede Univerze v Mariboru. Njeno habilitacijsko področje so kvantitativne metode v organizacijskih vedah. Ukvarja se z raziskavami na področju stohastičnih procesov, zanesljivosti in razpoložljivosti tehničnih sistemov ter varnosti informacijskih sistemov.