

VILNIAUS UNIVERSITETAS  
MATEMATIKOS IR INFORMATIKOS FAKULTETAS  
PROGRAMŲ SISTEMŲ KATEDRA

Baigiamasis bakalauro darbas

**Privačios informacijos išsaugojimas taikant dirbtinio intelekto  
technologijas**  
(Privacy-preserving AI)

Atliko: 4 kurso 3 grupės studentas

Paulius Milmantas (parašas)

Darbo vadovas:

dr. Linas Petkevičius (parašas)

Recenzentas:

(parašas)

Vilnius  
2020

## Turinys

Ivadas .....	2
1. Asmens duomenų privatumas mašininio mokymo kontekste .....	3
2. Pažeidžiamumo metrikos .....	6
2.1. Apibrėžimas .....	6
2.2. Skaičiavimas praktikoje.....	7
3. Modelio duomenų atgavimo metodai .....	8
4. Homomorfinis šifravimas .....	9
5. Federuotas mašininis mokymas .....	10
5.1. Metodas.....	10
5.2. Komunikacijos kaštų mažinimas .....	10
5.3. Įrenginių heterogeniškumas .....	10
5.4. Analizės heterogeniškumas .....	10
5.5. Saugos problemos .....	10
6. Blokų-grandinių technologija .....	11
7. Saugus skirtingų pusių skaičiavimas .....	12
8. Bendro įspūdžio privatus agregavimas .....	13

## Įvadas

Mašininis mokymas yra dirbtinio intelekto sritis, kuri pasitelkia statistinius algoritmus, kad apibrėžtų duomenų atpažinimo modelį. Modelio apmokymui dažniausiai reikia turėti daug duomenų. Kai kurie uždaviniai reikalauja duomenų, kurie nėra laisvai prieinami ir yra privatiūs. Mašininio mokymo tyrimų srityje yra kilusi problema dėl jų saugojimo. Vienas iš faktorių, kuris lėmė šį susidomėjimą yra 2016 metais Europoje priimtas duomenų apsaugos reglamentas (GDPR). Pagal jį, fizinių asmenų duomenys turi būti saugomi naudojantis tam tikromis taisyklėmis ir negali būti atskleisti trečiosioms šalims, be asmens sutikimo [1].

Šią problemą išspręsti siekia įvairūs tyrimai ir naujai atrasti metodai privatumą saugančio dirbtinio intelekto srityje. Šią problemą galima išskaidyti į kelias atskiras sritis:

- Analizuojamų duomenų privatumas [2]. Algoritmas apmoko modelį atpažinti duomenis. Turint sukurtą modelį, neturi būti galima atgaminti duomenų, pagal kuriuos jis buvo mokomas, bei negali būti identifikuoti asmenys. Taip nukentėtų žmonių privatumas ir būtų pažeistas Europos duomenų apsaugos reglamentas. Šio pažeidimo pavyzdys gali būti ir paprastas teksto atkūrimo modelis. Duodama sakinio pradžia, modelis nuspėja jo pabaigą. Jeigu suvedus tam tikras detales modelis užbaigia sakinį naudodamas asmeninius duomenis, kurie atskleidžia žmonių tapatybę, šis modelis nėra saugus.
- Duomenų įvesties privatumas [2]. Trečios šalys neturi matyti įvedamų duomenų. Tai gali būti tinklo saugumo spragos, duomenų surinkimo aplikacijų spragos ir t.t...
- Modelio išvesties privatumas [2]. Modelio išvesties neturi matyti asmenys, kuriems šie duomenys nepriklauso. Šis punktas yra sąlyginis, priklauso nuo modelio svarbos. Jeigu tai yra svarbūs asmeniniai duomenys, negalima rizikuoti. Tačiau jeigu tai yra viešai prieinami duomenys, šis punktas negalioja.
- Modelio apsauga [2]. Sukurtas modelis negali būti niekieno pasisavintas. Šis punktas yra skirtas apsaugoti programos kūrėją.

Darbo tikslas - ištirti ir palyginti privatumą saugančius dirbtinio intelekto algoritmus pagal jų saugumą, našumą ir panaudojamumą.

Darbo uždaviniai:

- Išanalizuoti esamus algoritmus pagal jų saugumą ir panaudojamumą.
- Ištirti kurie algoritmai yra realizuoti ir dalį nerealizuotų, realizuoti.
- Palyginti algoritmus pagal našumą.

## 1. Asmens duomenų privatumas mašininio mokymo kontekste

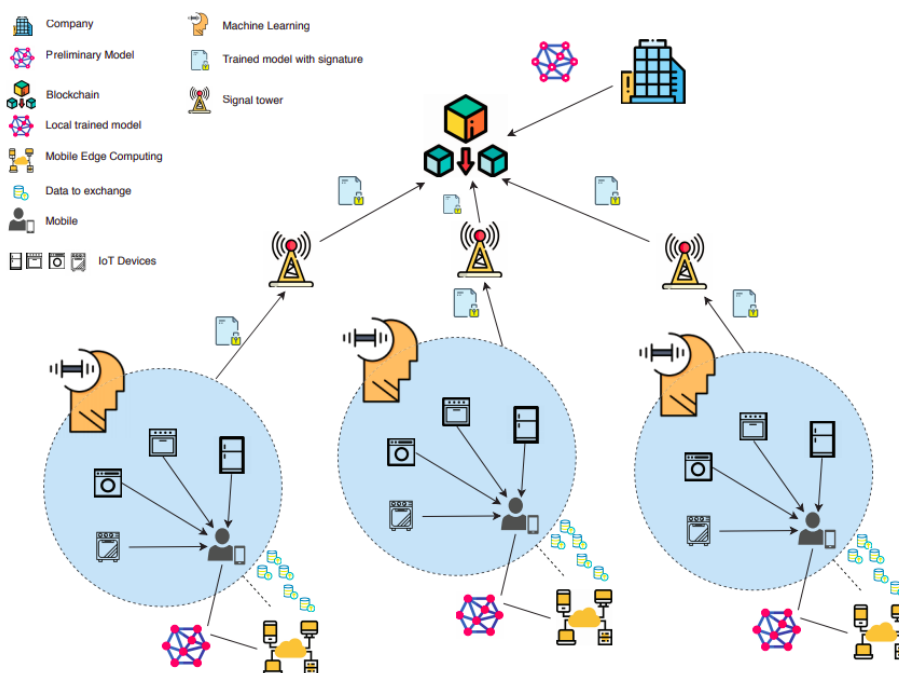
Vienas iš būdų apsaugoti duomenis yra decentralizuoti modelį ir naudoti paskirstyto mokymo algoritmus. Paprastai mašininiam mokymui yra naudojamas centralizuotas serveris. Yra surenkami duomenys į vieną vietą ir modelį moko vienas kompiuteris. Jeigu norima pridėti daugiau duomenų arba papildomą klasifikatorių, modelį reikia apmokyti iš naujo. Taip pat, kas kuria šį modelį, turi visą prieigą prie duomenų, šie modeliai nėra saugūs.

Šiai bėdai išspręsti, vienas iš metodų yra federuotas mašininis mokymas („Federated machine learning“). Šis metodas padeda išspręsti kaikurias bėdas, vykdant modelio mokymą decentralizuotai [3]. Prie bendro tinklo gali prisijungti daug įrenginių. Visi jie turi savo unikalų duomenų rinkinį. Kiekvienas įrenginys pasirenka geriausią statistinį metodą modelio mokymui ir jį sukuria. Taip yra sukuriami daug skirtingų modelių su skirtingais duomenų rinkiniais. Visi šie rinkiniai vėliau yra surenkami į vieną vietą ir toliau naudojami duomenų analizei. Taip surinkus atskirų įrenginių sukurtus modelius, neturime prieigos prie duomenų ir daugiau žmonių gali prisidėti prie modelio kūrimo, nematant pilno duomenų rinkinio. Tačiau naudojant šį metodą ne visos problemos yra išsprendžiamos. Duomenų saugumas nėra garantuojamas [3]. Jeigu duomenys nėra užšifruojami, jie gali būti pavogti. Tai gali būti padaryta, jeigu yra naudojamas nesaugus interneto ryšys, arba jeigu yra bandoma analizuoti atskirų įrenginių atsiųstus modelius. Šis metodas išsprendžia tik kelias problemas. Likusios yra: komunikavimo kaštai, įrenginių heterogeniškumas, statistinis heterogeniškumas ir saugumo problemos [4].

Visos šios problemos yra nagrinėjamos ir joms spręsti kuriami nauji metodai. Komunikacijos kaštams mažinti, yra kuriami algoritmai, kaip turi būti perduodamas modelis tinkle, kad būtų maža tinklo apkrova. Įrenginių heterogeniškumui spręsti, yra parenkami tam tikri įrenginių rodikliai ir pagal tai nustatoma, kokia bus įrenginio komunikacija: ar jis naudos asinchroninį komunikavimą, kaip dažnai tai vyks ir kokia yra įrenginio klaidos tikimybė [4]. Statistikos heterogeniškumui pataisyti kuriami yra metodai, kaip meta duomenų-mokymas ir keletos užduočių mokymas („Multi-task learning“) [4].

Straipsnyje apie privatumą išsaugančius blokų-grandinių tinklus („Privacy-Preserving Blockchain-Based Federated Learning for IoT Devices“) yra pateikiamas naujas paskirstyto mokymosi algoritmas, kuris yra paremtas blokų grandinių technologijomis [5]. Šis modelis yra grafiškai pateiktas pav.1. Šią architektūrą pavyzdyje sudaro 3 komponentai: gamintojai, klientai ir „block chain“ tinklas. Šiame pavyzdyje gamintojas pareiškia užklausa, kad reikia atlikti apklausą. Klientai, kurie sutinka su užklausa, išsiunčia savo sukurtą modelį. „Block chain“ tinklas elgiasi kaip centralizuotas serveris ir surenka visus klientų modelius. Tuomet, pasirinktas kompiuteris, kuris atlieka visą darbą („miner“) atlieka galutinį modelių sujungimą.

Minėti decentralizuoti būdai apjungia skirtingus sukurtus modelius į vieną ir visi įrenginiai kurie kuria modelius, naudoja savo unikalius duomenis. Saugus skirtingų pusių skaičiavimo metodas („Secure multiparty computation“) yra kriptografinis protokolas, kuris leidžia įrenginiams, su unikaliais duomenimis, skaičiuoti funkcijos reikšmę, nematant kitų įrenginių reikšmių [6]. Su šiuo



1 pav. Paskirstytas mokymas [5]

metodu, yra sukuriamas vienas modelis tarp įvairių įrenginių. Pagrindinis įrenginys, kuris ruošia skaičiavimo užklausą, gali suskaidyti pradinis duomenis, užšifruoti juos ir taip paskirstyti tarp kitų įrenginių. Taip sukurti įrenginio rezultatai negali būti atšifruoti ir pasisavinti, o minėto federuoto mokymo metodo sukurti rezultatai, jeigu nėra gerai apsaugoti, gali būti pasisavinti ir atgauti pradiniai duomenys su kuriais modelis buvo sukurtas.

Kitas būdas apsaugoti modelį yra naudoti homomorfinį šifravimą. Pagal šį metodą, modelis yra mokomas naudojant užšifruotais duomenimis. Jie mokymo metu, nėra atšifruojami. Šis metodas leidžia keliems įrenginiams vienu metu atlikti skaičiavimus su duomenimis, kurių jie nemato. Yra daug metodo varijacijų. Kaikurios yra pažeidžiamos ir gali atskleisti visos infrastruktūros duomenis [7]. Dažniausiai homomorfinės sistemos yra labiau pažeidžiamos nei nehomomorfinės [7].

2016 metais buvo pasiūlytas naujas metodas paspartinti pilnai homomorfinės šifravimo sistemos [8]. Homomorfinės sistemos, kurios naudoja šį metodą yra laikomos ketvirtos kartos. Šis metodas leidžia aproksimuoti užšifruoto teksto sudėtį, daugybą ir pakeisti atšifruoto teksto proporcijas [8]. Proporcijų pakeitimo procedūra suskaido užšifruotą tekstą į dalis, to pasekoje yra apvalinamas neužšifruotas teksas. Šio metodo pagrindinė idėja yra pridėti triukšmo prie pagrindinės žinutės. Šis triukšmas yra pridodamas prie neužšifruoto teksto dėl saugumo priežasčių ir jis yra laikomas kaip aproksimavimo paklaida. Tuo pasekoje, metodas veikia su tam tikra paklaida

Bendro išspūdžio agregavimo metodas yra naujausias iš visų minėtų. Jis leidžia įrenginiui su savo duomenimis sukurti modelį naudojant betkokių metodus ir taip prisidėti prie bendro progreso [6]. Visi modeliai yra surenkami ir leidžiami duomenys per šiuos modelius. Kiekvienas modelis skiria savo balsą, jie surenkami ir padaromas bendras sprendimas. Šis metodas leidžia

lengvai plėsti modelį. Taip pat yra saugomas ir duomenų privatumas. Jeigu keli modeliai, kurie nesidalina duomenimis, teigia vienodai, tai reiškia, kad negalima atkurti duomenų ir jie yra saugūs [6]. Skaičiavimas vykdomas pasirenkant optimalią strategiją ir imant mažiausią galimą žingsnių skaičių.

## 2. Pažeidžiamumo metrikos

### 2.1. Apibrėžimas

Norint apsaugoti duomenis, sukurtam duomenų aptikimo modeliui reikia atlikti analizę, kaip tikėtina, kad modelio pradiniai duomenys bus atkurti. Tai išanalizuoti yra daug būdų, vienas iš jų yra pateiktas straipsnyje [9]. Šis straipsnis ieško metrikos reikšmės, kuri parodo, kaip pradinis duomenis atsimena modelis.

Tarkime turime duomenų rinkinį  $s[r]$ . Šiam rinkiniui pirma reikia apskaičiuoti rangą [8].

$$rangas(s[r]) = |\{r' \in R : Px(s[r']) \leq Px(s[r])\}| \quad (1)$$

Šis rangas nurodo, kurioje vietoje yra šis inicijuotas duomenų rinkinys sąraše, tarp visų galimų testinių rinkinių kombinacijų. Pavyzdžiui, norima sukurti natūralios kalbos atpažinimo modelį. Jo mokymui, yra paduodamas atsitiktinai sugeneruotas sakinyss „Atsitiktinis skaičius 125“ ir apskaičiuojama jo entropija. Tada galima išrašyti visus galimus sakinius, kuriuos yra leidžiama siųsti, ir juos surūšiuoti pagal entropijos laipsnį. Tada reikia apskaičiuoti, kaip dažnai pasitaiko panašūs sakiniai. Taip yra gaunamas rangas. Kadangi šis pavyzdis sakinyss yra sudėtingas, su atsitiktiniais skaičiais ir jo entropija yra aukšta, tikriausiai aukštesnė nei dauguma mokymo duomenų, jo rangas bus artimas vienetui. Darant prielaidą, kad šio sakinio entropija yra pati aukščiausia iš galimų testavimo duomenų, galima teigti, jog šio sakinio rangas yra lygus vienetui.

Rangas tiesiogiai nepasako kokia yra tikimybė, kad bus sugeneruotas toks testavimo rinkinys. Jo skaičiavimas paima daug resursų, nes reikia sugeneruoti visas galimas kombinacijas [8].

Kitas nagrinėjamas kintamasis yra „atvirumo“ metrika. Tai aproksimuotas skaičius, kuris nurodo, kaip tikėtina, jog testavimo duomenys bus atgaminti iš modelio. Ši metrika pasako ką naujo sužinome apie pradinis duomenis, kai per modelį paleidžiame atsitiktinai sugeneruotus duomenis. Todėl, reikia skaičiuoti spėjimo entropijos redukcijos lygtį. Spėjimo entropija, tai spėjimų skaičius  $E(X)$ , kuris reikalingas atspėti diskretų, atsitiktinai sugalvotą parametrą  $X$ .

Jeigu kintamasis  $r$  yra pasirenkamas atsitiktinai  $r \in R$ , tai reiškia, kad reikia generuoti atsitiktinius skaičius tol, kol bus gauta  $r$  reikšmė. Iš to seka, kad spėjimų turėtų būti lygus [2] lygčiai.

$$E(s[r]) = \frac{1}{2}|R| \quad (2)$$

Turint suskaičiuotus galimus rinkinių rangus, bei entropijas, galima naudoti patobulinta skaičiavimo taktiką. Visi galimi duomenys yra surūšiuojami pagal entropiją arba sudėtingumą. Sąrašo pradžioje yra elementas su mažiausia entropija. Jo atspėjimo tikimybė yra viena iš didžiausių. To pasekoje, gauname formulę [3].

$$E(s[r]|f_\theta) = rank_\theta(s[r]) \quad (3)$$

Ši formulė padidina skaičiavimų spartumą ir tai galima apskaičiuoti padalinus vieną formulę

iš kitos. Tai pateiktas [4] lygtyje.

$$\frac{E(s[r])}{E[s[r]|f_\theta]} = \frac{\frac{1}{2}|R|}{rank_\theta(s[r])} \quad (4)$$

Rangai tiksliai neapibrėžia kokia yra tikimybė, jog elementas bus atspėtas. Jis tik lygina visas galimas kombinacijas ir skaičiuoja jų entropijas. To pasekoje, visi šie skaičiavimai yra tik apytikslis spėjimas. Kadangi tai nėra tikslu ir norima sužinoti tik bendrą vaizdą apie algoritmą, galima naudoti logaritmus. Tai atlikta lygtyje [5].

$$\log_2\left(\frac{E(s[r])}{E[s[r]|f_\theta]}\right) = \log_2\left(\frac{\frac{1}{2}|R|}{rank_\theta(s[r])}\right) \quad (5)$$

Šias lygtis supaprastinus, gaunama lygtis [6].

$$\log_2|R| - \log_2(rank_\theta(s[r])) - 1 \quad (6)$$

Dėl paprastumo yra pridamas vienetas ir taip gaunama atvirumo metrika, pateikta lygtyje [7].

$$atvirumas(s[r])_\theta = \log_2|R| - \log_2 rangas_\theta(s[r])$$

(7)

## 2.2. Skaičiavimas praktikoje

Pažeidžiamumo metrika „Atvirumas“ yra naudojama tikrinant mašininio mokymo algoritmus duomenų įsiminimą. Atlikus šiuos skaičiavimus, programuotojai gali spręsti dėl tolimesnių veiksmų: tęsti kūrimą ar dirbti ties duomenų anonimizavimu.



### **3. Modelio duomenų atgavimo metodai**

## **4. Homomorfinis šifravimas**

## **5. Federuotas mašininis mokymas**

### **5.1. Metodas**

### **5.2. Komunikacijos kaštų mažinimas**

### **5.3. Įrenginių heterogeniškumas**

### **5.4. Analizės heterogeniškumas**

### **5.5. Saugos problemos**

## **6. Blokų-grandinių technologija**

## **7. Saugus skirtingų pusių skaičiavimas**

## **8. Bendro įspūdžio privatus agregavimas**