

VILNIAUS UNIVERSITETAS
INFORMATIKOS INSTITUTAS
PROGRAMŲ SISTEMŲ KATEDRA

**Anomalių aptikimas interneto prieigos
stebėsenos sistemos įrašuose**

**Detection of anomalies in internet access monitoring
system data**

Bakalauro darbas

Atliko:	Jokūbas Rusakevičius	(parašas)
Darbo vadovas:	asist. dr. Vytautas Valaitis	(parašas)
Darbo recenzentas:	j. asist. Linas Petkevičius	(parašas)

Vilnius – 2019

Santrauka

TODO: Santrauka

Raktiniai žodžiai: „MacroBase“, IPSS, analitinis įrankis, raktinis žodis 4, raktinis žodis

5

Summary

TODO: summary

Keywords: „MacroBase“, IPSS, analitic engine, keyword 4, keyword 5

TURINYS

IVADAS	6
Problematika	6
0.1. Darbo tikslai ir uždaviniai	7
1. „MACROBASE“ ANALITINIS ĮRANKIS	9
1.1. Apie „MacroBase“	9
1.2. „MacroBase“ sprendžiamos problemos	10
2. „MACROBASE“ „SQL“ PAKETO VEIKIMO PRINCIPAS – „MACROBASE SQL“	11
3. DUOMENŲ ANALIZĖS EKSPERIMENTAS NAUDOJANT „MACROBASE SQL“ MODULĮ	12
3.1. Duomenų rinkinys	12
3.1.1. Eksperimentui pritaikyti duomenys	12
3.1.1.1. Duomenų aprašymas	12
3.1.1.2. Pradinių duomenų paruošimas	15
3.1.1.3. Papildytų duomenų paruošimas: koordinacių suapvalinimas, aukščio virš jūros lygio pridėjimas	16
3.1.1.4. Papildytų duomenų paruošimas: metų ir valandų išskyrimas	17
3.1.1.5. Papildytų duomenų paruošimas: apvalinimo variacijos	18
3.2. Eksperimento aplinkos paruošimas	19
3.2.1. Virtuali mašina eksperimentui	19
3.2.1.1. Virtualios mašinos paruošimas	20
3.2.2. „MacroBase“ analitinio įrankio paruošimas	20
3.2.2.1. „MacroBase SQL“ diegimui reikalingi papildomi paketai	21
3.2.2.2. „MacroBase SQL“ analitinio įrankio diegimas	21
3.3. Eksperimento vykdymo eiga	22
3.3.1. Konfigūracija	22
3.3.2. Eksperimento eiga	22
3.3.2.1. Eksperimento eiga su pritaikytais IPSS belaidžio interneto 3G duomenimis	23
3.3.2.2. Eksperimento eiga su pritaikytais IPSS belaidžio interneto 3G duomenimis, papildytais aukščio virš jūros lygio ir suapvalintų reikšmių laukais	23
3.3.2.3. Eksperimento eiga su pritaikytais IPSS belaidžio interneto 3G duomenimis, papildytais metų ir valandų laukais	24
3.3.2.4. Eksperimento eiga su pritaikytais IPSS belaidžio interneto 3G duomenimis, tikrinami optimaliausi reikšmių apvalinimai	24
3.3.2.5. Eksperimento eiga su pritaikytais IPSS belaidžio interneto LTE, pritaikius išanalizuotus 3G parametrus	25
3.3.2.6. Eksperimento eiga su pritaikytais IPSS belaidžio interneto WiMAX, pritaikius išanalizuotus 3G parametrus	25
3.3.3. Eksperimentinio tyrimo rezultatų analize	25
REZULTATAI IR IŠVADOS	26
LITERATŪRA	27
SANTRUMPOS	29
PRIEDAI	29
1 priedas. IPSS matavimų koordinatės 2014–2018 metais	30
2 priedas. Pradinių IPSS CSV failų apdorojimo „Python“ skriptas	37

3 priedas. Atrinktų IPSS CSV failų duomenų papildymo aukščio virš jūros lygio ir suap- valintų reikšmių laukais „Python“ skriptas	39
4 priedas. Atrinktų IPSS CSV failų duomenų papildymo metų ir valandų laukais „Python“ skriptas	43
5 priedas. Atrinktų IPSS papildomų CSV failų generavimo „Python“ skriptas.....	45
6 priedas. „MacroBase SQL“ SQL užklausų pavyzdys.....	47

Įvadas

Šis darbas yra Programų Sistemų studijų baigiamasis bakalauro darbas apie anomalijų aptikimą interneto prieigos stebėsenos sistemos įrašuose naudojant „MacroBase“ atviro kodo analitinį įrankį bei „SQL“ paketą.

Problematika

Informacijos rinkimas fizine forma, juos užrašant ant popieriaus lapo, pasak MIT Media Lab direktoriaus Joi Ito, dar visai neseniai buvo vienintelis būdas rinkti, klasifikuoti ir saugoti duomenis. Žmogus mintis ir idėjas užrašydamas ant popieriaus lapo – paversdavo žiniomis (angl. *knowledge*). Tačiau laikai keičiasi, ir dabartinė didžiųjų duomenų (angl. *big data*) situacija yra kitokia. Dėl milžiniškų surenkamų duomenų kiekio, nėra trivialu atskirti naudingus ir nenaudingus duomenis, todėl, natūraliai, negalima visų surenkamų duomenų automatiškai priskirti prie ir vadinti žiniomis. Surinkti duomenys nėra žinios, kol jų nepradedama nagrinėti, analizuoti ir filtruoti, ir tik pradėjus šį procesą galima pastebėti, kad gaunama informacija yra įdomi ir netgi svarbi [Smo14].

Duomenų, kuriais yra operuojama, juos yra renkama ir saugoma, kiekiai nuolatos didėja [EMC14]. Iš tiesų, dėl daugelio priežasčių net greitis, kuriuo šie duomenys yra generuojami eksponentiškai kyla [Mak17]. Jau 2015–2016 metais didžiosios socialinių tinklų kompanijos Twitter, Facebook ir LinkedIn pranešė fiksuojančios iki 12 milijonų įvykių (angl. *events*) per sekundę [Ast16; PFT⁺15; Woo15]. Šie duomenų kiekiai gerokai lenkia žmogaus gebėjimą juos apdoroti bei žymiai apsunkina darbą tiek analitikams, tiek analitiniams įrankiams. Net aukščiausios kvalifikacijos (angl. *best-of-class*) šių dienų taikomųjų programų operatoriai bei analitikai praneša panaudojantys anekdotiškai mažą jų surenkamų duomenų kiekį – mažiau nei 6% [BGR⁺17b].

Renkant tokius didelius kiekius duomenų natūraliai, užfiksuojami įrašai, kurie yra nebūdingi įrašų rinkiniui bei išsiskiria. Duomenų vienetai nukrypę nuo kitų duomenų rinkinyje yra vadinami anomalijomis. Anomalijų aptikimas yra procesas, kurio metu yra aptinkama ir identifikuojama anomalija arba išskirtis (angl. *outlier*) duomenų rinkinyje [Tec19]. Šio darbo metu bus tiriamas anomalijų aptikimas remiantis belaidės interneto prieigos duomenų perdavimo spartos kontrolinių matavimų, atliktų 2014–2018 metais, rezultatais, kurie gauti naudojantis Lietuvos Respublikos ryšių reguliavimo tarnybos administruojama IPSS [rrtar].

Fiksuojami esminiai ar sisteminiai duomenys tobulu atveju turėtų būti peržiūrimi reguliariai ir pilnai, tam kad būtų galima užtikrinti tolygią sistemų veiklą, matavimų taisyklingumą ar renkamų duomenų korektiškumą ir laiku identifikuoti bei eliminuoti kylančias klaidas bei užkirsti

kelią kenkėjiškoms atakoms. Tačiau kaip jau minėta anksčiau, didėjantis duomenų kiekis gerokai lenkia žmogaus gebėjimą juos apdoroti. Duomenų kiekis nuolatos didėja, tačiau žmogaus dėmesys išlieka riboto dydžio. Iškyla iššūkis prioritetizuoti žmogaus dėmesį – akcentuoti, filtruoti, jungti, grupuoti ir kontekstualizuoti svarbiausius analizuojamų duomenų elgsenos ir būsenos pokyčius, taip padedant išskirti svarbiausius duomenų ypatumus ir rodyti žmogui tik apibendrintą, svarbia ir ribotą kiekį informacijos [BGR⁺17b]. Visa nereikalinga informacija, kuri yra rodoma žmogui reikalauja ir eikvoja jo dėmesį [Sim71]. Nors žmonėms yra fiziškai neįmanoma peržiūrėti visų duomenų, mašinos ir/ar kompiuteriai gali tai atlikti, ar bent gebėti, jei jų resursai būtų ekonomiškai paskirstyti. Analitikų teigimu, didelis kiekis svarbių duomenų yra praleidžiama dėl neefektyvaus ir/ar lėto mašinų darbo [BGR⁺17b].

2017 metais Standfordo universitetas kartu su Masačusetso technologijos institutu paskelbė kuriantys naują atviro kodo analitinį paieškos įrankį „MacroBase“. Šis įrankis yra skirtas ne tik darbui su didžiaisiais bet ir su greitaisiais duomenimis (angl. *fast data*). Pagrindiniu uždaviniu „MacroBase“ sau išsikėlė žmogaus dėmesio prioritetizavimą. Vienas iš daugelio „MacroBase“ šio uždavinio sprendimų yra didžiausio dėmesio reikalaujančios išvesties sugeneravimas ir jos supaprastinimas, kur neįprastus duomenų vienetus „MacroBase“ padeda aiškinti pagal duomenų atributus [BGM⁺17; BGR⁺17a; BGR⁺17b]. „MacroBase“ yra pakankamai naujas analitinis įrankis, gebantis grąžinti tikslius rezultatus dirbdamas 2 milijonų įvykių per sekundę, per užklausą, per branduolį greičiu, kuris taip pat pataikia daug inovatoriškų sprendimų vis didėjančių ir greitėjančių duomenų srautų analizei atlikti [BGM⁺17]. Šiame darbe bus plačiau nagrinėjamas bei eksperimentai atliekami su „MacroBase“ įrankiu. Plačiau apie „MacroBase“ 1 ir 2 skyriuose.

0.1. Darbo tikslai ir uždaviniai

Šio darbo **tikslas** – ištirti „MacroBase“ duomenų analizavimo ir anomalijų aptikimo įrankio „SQL“ paketą – „MacroBase SQL“ ir pritaikyti jo veikimą analizavimui ir anomalijų aptikimui IPSS belaidės interneto prieigos duomenų perdavimo spartos kontrolinių matavimų įrašuose.

Darbui išskelti **uždaviniai**:

1. Išanalizuoti „MacroBase SQL“ veikimo principą.
2. Paruošti eksperimentinę aplinką ir IPSS įrašų rinkinį eksperimentui.
3. Palyginti eksperimentų rezultatus, keičiant pagrindinių parametrų (procentilio, palaikymo ir rizikos koeficiento) reikšmes.
4. Atlikti duomenų transformacijas (pridėti, atimti, filtruoti, keisti reikšmes) ir palyginti eksperimento rezultatus.

5. Atlikti tas pačias transformacijas su kitais panašiais duomenimis ir palyginti eksperimento rezultatus.
6. Pateikti rekomendacijas IPSS duomenų anomalijų aptikimui naudojantis „MacroBase“.

1. „MacroBase“ analitinis įrankis

Šiame skyriuje yra aprašomas „MacroBase“ analitinis įrankis (1.1 poskyrius). Įrankis kuriamas Standfordo Universiteto bedradarbiaujant kartu su Masačusetso technologijos institutu. Taip pat, šiame skyriuje rašoma apie „MacroBase“ žmogaus dėmesio prioritetizavimo bei kitų problemų sprendimo būdus (1.2 poskyrius).

1.1. Apie „MacroBase“

„MacroBase“ nuo pradžių buvo kurtas būti žmogaus dėmesio prioritetizavimui skirtas analitinis įrankis, bet kad tai būtų pasiekta, reikia naudoti efektyvius aiškinamuosius ir klasifikavimo statistinius operatorius, ir dirbti su duomenų srautais [Sur17]. „MacroBase“ naudoja du pagrindinius operatorių tipus:

- **Klasifikavimas** – tai operatorius, kuris pagal naudotojo nurodytas klases pats susižymi individualius taškus, kuriuos jau yra ištyręs (pavyzdžiui, statistiškai „normalus“ arba statistiškai „išskirtinis“).
- **Aiškinamieji** – tai operatoriai, kurie grupuoja ir jungia duomenų taškus po kelis.
Toliau trumpai aprašoma „MacroBase“ darbo proceso struktūra [BGM⁺17; Rus18]:
 - Duomenų gavimas (angl. *ingestion*):
 - „MacroBase“ gauna išorinį duomenų srautą analizei.
 - Duomenų taškai susideda iš metrikų – matavimų (pavyzdžiui, greitis, kelionės laikas, baterijos procentai), ir atributų – siejamų duomenų (pavyzdžiui, telefono numeris, kompiuterio MAC adresas, paso kodas).
 - „MacroBase“ anomalijų aptikimui naudoja metrikas, kurias po to paaiškina naudodama atributus.
 - Duomenų transformavimas (angl. *feature transformation*):
 - Tai yra neprivalomas procesas.
 - Šio proceso metu „MacroBase“ vykdo domenui specifinių duomenų transformacijas. Šios transformacijos gali būti įvairių tipų: statistinės, priklausančios nuo duomenų tipo ar kt.
 - Duomenų klasifikavimas (angl. *classification*):
 - Kiekvienas duomenų taškas yra klasifikuojamas.
 - „MacroBase“ šiuos taškus klasifikuoja jiems priskirdama po žymę, kuri yra priklausomą nuo metrikos.
 - Taip pat, šioje proceso stadijoje yra treniruojamai klasifikatoriai. Klasifikatoriai yra

treniruojami ir vertinami pagal srauto metrikas.

- Duomenų aiškinimas (angl. *explanation*):
 - „MacroBase“ negrąžina tiesioginio sužymėtų duomenų taškų.
 - „MacroBase“ sugeneruoja paaiškinimus ir apjungia sužymėto srauto duomenų taškus.
 - Kadangi duomenų srauto duomenys yra nuolatos apdorojami, „MacroBase“ aiškinamieji operatoriai grąžina paaiškinimus tik gavusi naudotojo užklausą, taip išvengdamas nereikalingų skaičiavimų.
- Duomenų pateikimas (angl. *presentation*):
 - Statistiniai reitingai yra priskiriami, dėl galimai didelio aiškinimų kiekio.
 - „MacroBase“ naudotojams aiškinimus pateikia jų svarbos tvarka.

1.2. „MacroBase“ sprendžiamos problemos

2. „MacroBase“ „SQL“ paketo veikimo principas – „Macro-Base SQL“

3. Duomenų analizės eksperimentas naudojant „MacroBase SQL“ modulį

Eksperimentui atlikti buvo pasirinktas „MacroBase SQL“ analitinis įrankis „MacroBase“ su „SQL“ moduliu. Kadangi „MacroBase SQL“ yra leidžiama per komandinę eilutę, grafinė sąsaja kaip „MacroBase GUI“ naudota nebuvo. Atliekant eksperimentą buvo pasirinkti CSV tipo failų duomenų šaltiniai.

3.1. Duomenų rinkinys

Šiame skyriuje aprašyti eksperimentui pritaikyti duomenų rinkiniai.

3.1.1. Eksperimentui pritaikyti duomenys

Eksperimentui buvo naudojami Lietuvos Respublikos ryšių reguliavimo tarnybos administruojamos IPSS atliekami belaidės interneto prieigos duomenų perdavimo spartos kontrolinių matavimų 2014–2018 metų rezultatai. Matavimai atlikti operatorių UAB „Bitė Lietuva“, AB „Telia Lietuva“, UAB „TELE2“ ir AB LRTC judriojo ryšio tinkluose visoje Lietuvos teritorijoje. Matavimų paskirtis yra stebėti ir įvertinti teikiamų interneto prieigos paslaugų kokybę, kaip to reikalauja Europos Sąjungos direktyvos, taip pat supažindinti visuomenę su tokių matavimų rezultatais. Įrašai pateikiami CSV formato failais. Iš viso gauti 3 failai skirti trimis skirtingoms technologijoms: 3G, LTE ir WiMAX. Viena CSV faile priklausomai nuo failui priskirtos technologijos buvo apytiksliai 136000, 157000 ir 18000 įrašų eilučių, vidutiniškai 104000. Vieno failo dydis vidutiniškai 8,5MB.

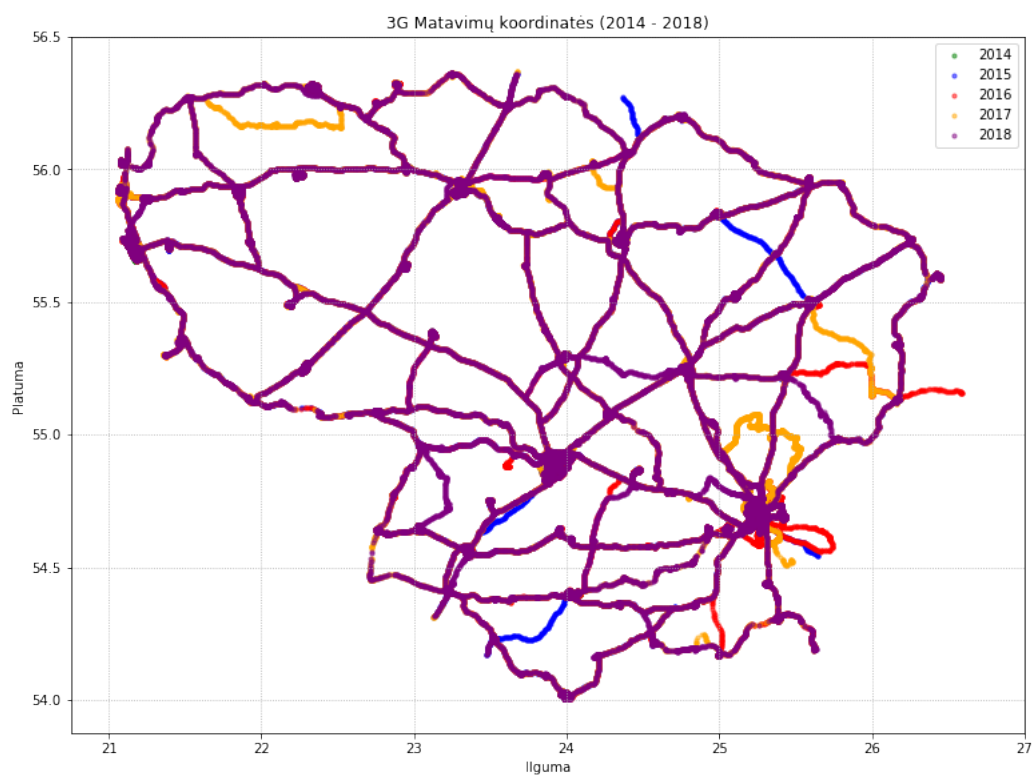
3.1.1.1. Duomenų aprašymas

Kiekvienas IPSS įrašų failas prasideda antraštės (angl. *header*) eilute, kuri susideda iš faile saugomų įrašų sąrašo laukų (angl. *fields*) pavadinimų. IPSS įrašų sąrašų laukai ir jų paaiškinimas:

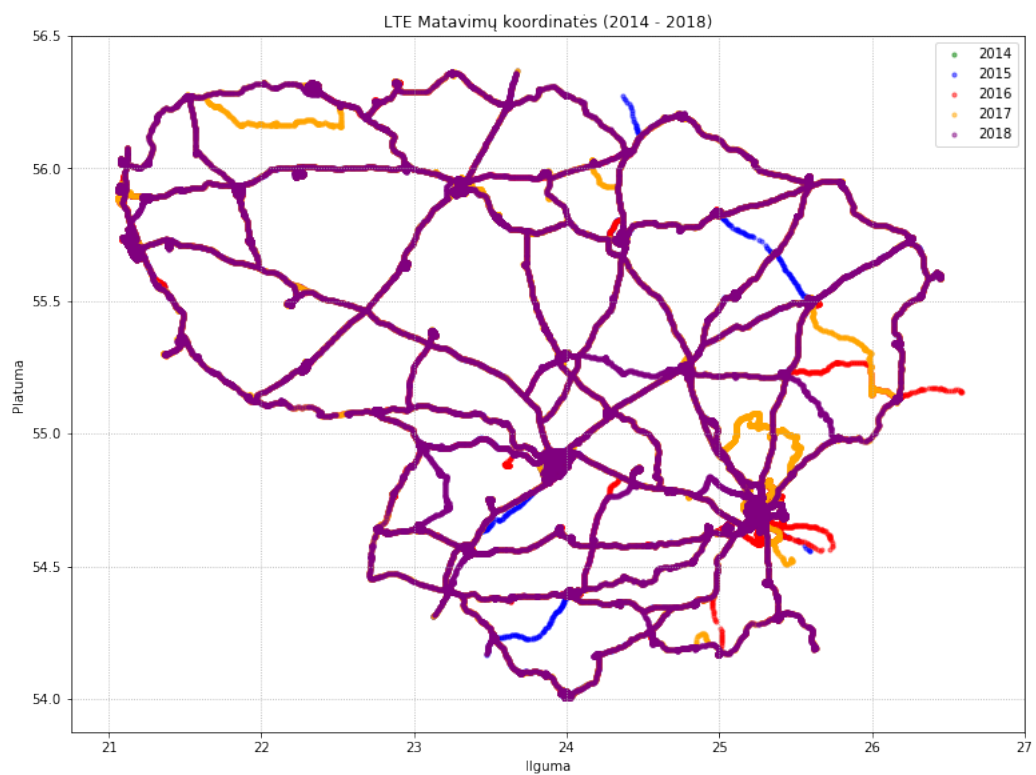
- **Data ir laikas** – įrašo fiksavimo data ir laikas, užrašoma „yyyy-MM-dd hh:mm:ss“ formata.
- **Platuma** – matavimo įrašo fiksavimo koordinatijų platuma užrašoma dešimtainiais laipsniais (angl. *decimal degrees*).
- **Ilguma** – matavimo įrašo fiksavimo koordinatijų ilguma užrašoma dešimtainiais laipsniais (angl. *decimal degrees*).

- **Operatorius** – matavimo įrašo operatoriaus skaitinis kodas arba pavadinimas (vieno iš jau minėtų operatorių: „Bitė Lietuva“, „Telia Lietuva“, „TELE2“ ar LRTC).
- **Ryšio tinklas** – tinklo technologija (3G, LTE arba WiMAX).
- **Celės id** – korinio tinklo ląstelės identifikacinis numeris atvaizduotas skaitine šešiolyktaine forma.
- **Bazinės stoties id** – bazinės stoties identifikacinis numeris atvaizduotas skaitine šešiolyktaine forma, parametras pateikiamas tik LRTC WiMAX tinkle atliktuose matavimuose.
- **3G ryšio technologija** – matavimų metu nustatytas duomenų perdavimo technologijos porūšis (WCDMA, HSPA+, DC-HSPA+).
- **RSSI** – (angl. *Received Signal Strength Indicator*) signalo lygis, dBm (apklausiant modemą gaunamos vertės nuo -51 dBm iki -113 dBm. Mažiausia signalo vertė -113 dBm reiškia labai silpną signalą, kuriam esant praktiškai nebegalima naudotis duomenų perdavimo paslauga).
- **RSRP** – (angl. *Reference Signal Received Power*) pilotinio signalo galia, dBm (apklausiant modemą gaunamos vertės nuo -140 dBm iki -44 dBm. Mažiausia signalo vertė -140 dBm reiškia labai silpną signalą, kuriam esant praktiškai nebegalima naudotis duomenų perdavimo paslauga). Parametras pateikiamas tik LTE tinkle atliktuose matavimuose.
- **RSRQ** – (angl. *Reference Signal Received Quality*) pilotinio signalo kokybė, dB (apklausiant modemą gaunamos vertės nuo -19,5 dB iki -3 dB. Mažiausia vertė -19,5 dB). Parametras pateikiamas tik LTE tinkle atliktuose matavimuose.
- **SINR** – (angl. *Signal to Interference and Noise Ratio*) signalo ir trukdžių plus triukšmo santykis, dB (apklausiant modemą gaunamos vertės nuo -20 dB iki 30 dB. Mažiausia vertė -20 dB). Parametras pateikiamas tik LTE tinkle atliktuose matavimuose.
- **CINR** – (angl. *Carrier to Interference and Noise Ratio*) nešlio ir trukdžių plus triukšmo santykis, dB, parametras pateikiamas tik LRTC WiMAX tinkle atliktuose matavimuose.
- **Sparta kbit/s** – (angl. *Download speed*) išmatuota duomenų gavimo spartos vertė, kbps.

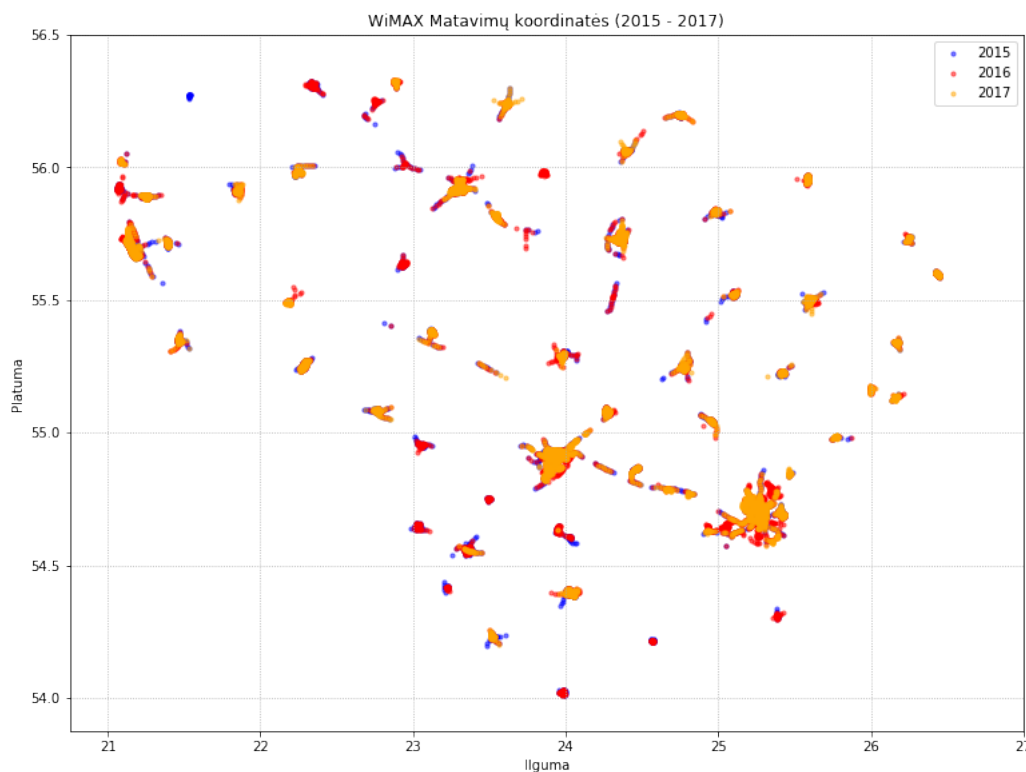
IPSS matavimai atliekami naudojant Ryšių reguliavimo tarnybos turimą įrangą, kuri yra sumontuota matavimams skirtame automobilyje. Matavimai atliekami automobiliui judant (angl. *drive test*) miestų gatvėmis, automagistralėmis arba rajoniniais keliais pagal pasirinktus maršrutus. Atliktų matavimų maršrutų bei koordinatės galima pamatyti 1, 2 ir 3 paveikslėliuose.



1 pav. 3G Matavimai 2014–2018 metais



2 pav. LTE Matavimai 2014–2018 metais



3 pav. WiMAX Matavimai 2015–2017 metais

Aiškiau išskirtus įrašų koordinacių duomenis galima rasti Priede nr. 1.

3.1.1.2. Pradinių duomenų paruošimas

Eksperimentui atlikti buvo atrinkti laukai, kurie turi reikšmes visose įrašų sąrašo eilutėse bei taip pat buvo atmesti pasikartojančios informacijos laukai (ryšio tinklo ir ryšio technologijos laukai, kurie nurodo, koks tinklo tipas iš trijų tiriamųjų yra naudojamas faile). Taigi, 3G tinklui atrinkti šie pradiniai laukai: data, laikas, ilguma, platuma, operatorius, celės id, rssi, 3G ryšio technologija ir sparta; iš LTE duomenų atrinkti laukai: data, laikas, ilguma, platuma, operatorius, celės id, rssi, sparta; iš WiMAX duomenų atrinkti laukai: data, laikas, ilguma, platuma, operatorius, bazinės stoties id, rssi, sparta.

Eksperimentui atlikti buvo reikalingas CSV formato failai, kuriuose nebūtų lietuviškų raidžių ir būtų naudojami vieno žodžio arba vienos simbolių eilutės antraštės pavadinimai (pavyzdžiui, tarpai gali būti pakeičiami „_“ simboliu), todėl buvo realizuotas „Python“ skriptas pavadintas „RenameAndRemoveFields.py“, kuris nuskaito duotųjų duomenų CSV failus ir juos apdorojus, grąžina naujus CSV failus. Skriptas, kaip ir visi darbo metu kurti skriptai, kurtas ir leistas „Jupyter Notebook“ aplinkoje. Skripte galima pasirinkti kelį iki duotųjų CSV failų ir kelį iki naujai sukuriamų CSV failų. Taip pat galima nurodyti norimas antraštes ir norimus laukus. Šio skripto paskirtis yra pervadinti antraštes ir atmesti nereikalingus laukus ir atskirti datą nuo

laiko. Skriptas atlieka šiuos veiksmus 3 kartus kiekvienai ryšio technologijai (3G, LTE, WiMAX) (Priedas nr. 2):

1. Nuskaito CSV įrašų failą pateiktą nurodytu keliu (18–23 ir 38 eilutės).
2. Sukuria naują CSV failą nurodytų keliu ir pavadinimu (24 ir 39 eilutės).
3. Sukuria „csv.writer“ objektą skirtą duomenų rašymui į CSV failus (25 eilutė).
4. Įrašo į naujai sukurtą CSV failą, naują antraštę, su naujais angliškais vieno žodžio (simbolių eilutės) pavadinimais (26 eilutė).
5. Vykdo ciklą, kuris iš pradinio failo nuskaitytų duomenų kiekvieno įrašo atrenka reikalingus laukus, atskiria laiką nuo datos ir įrašo į CSV failą naujai sukurtą eilutę. Taip pat atlieką kiekvienos eilutės skaičiavimą (28–34 eilutės).
6. Išvedamas galutinis eilučių skaičius (35 eilutė).
7. Į konsolę išvedama operacijos užbaigimo žinutė (40 eilutė).

3.1.1.3. Papildytų duomenų paruošimas: koordinačių suapvalinimas, aukščio virš jūros lygio pridėjimas

IPSS duomenyse yra pateikta daug informacijos, tačiau, dėl didelės jų įvairovės ir dėl realių skaičių laukuose, tik maža dalis jų yra potencialiai panaudojama aiškinimo procesui. Todėl buvo sukurtas antras skriptas pradinių duomenų papildymui (papunktis 3.1.1.2). Skripto paskirtis buvo pridėti papildomus koordinačių laukus, kuriuose būtų tos pačios koordinatės tačiau suapvalintos iki pasirinkto skaičiaus po kablelio. Skriptas turėjo ir antrą uždavinį – pasinaudojus tiksliais, originaliomis koordinatėmis rasti aukštį virš jūros lygio ir pridėti šią bei suapvalintą reikšmę į joms skirtus naujus laukelius. Šis skriptas buvo pavadintas „AddAltitudeAndRoundedColumns.py“. Skriptas leidžia pasirinkti kelią iki pradinio CSV failo, kelią iki naujai sukuriama CSV failo, taip pat skaitmenų po kablelio apvalinimo skaičių, bei aukščio virš jūros lygio apvalinimo sąlygas. Kadangi skriptas naudoja „Jawg.io“ API, jam reikia pateikti vartotojo prieigos žymę (angl. *access token*). Taigi naujas „Python“ skriptas nuskaito prieš tai sukurtus CSV failus, gauna aukščio virš jūros lygio reikšmes, suapvalina koordinačių bei aukščio reikšmes ir sukuria naujus CSV failus pasirinktu keliu ir pavadinimu. Skriptas atlieka šiuos veiksmus 3 kartus kiekvienai ryšio technologijai (3G, LTE, WiMAX) (Priedas nr. 3):

1. Nuskaito CSV įrašų failą pateiktą nurodytu keliu bei atskiria antraštę nuo įrašų (9–13 ir 89 eilutės).
2. Iš antraštės surandami platumos ir ilgumos laukų indeksai (14–17 ir 90 eilutės).
3. Pasinaudojus gautais indeksais, sukuriama visų ilgumos ir platumos reikšmių masyvai (18–21 ir 91 eilutės).

4. Kadangi bus atliekamos HTTP užklauso ir URL užklauso ilgis yra ribots, yra vykdomas ciklas, kuris visus įrašus padalina į pasirinkto žingsnio dydžio dalis (22–34 ir 92).
5. Vykdomas ciklas cikle, kuris visas ilgumos ir platumos reikšmes sujungia į „string“ tipo kintamuosius, kuriuose ilguma ir platuma yra atskirtos kableliu „|“ (26–31 eilutės).
6. Kiekvienam žingsniui, iš platumos ir ilgumos bendrų kintamųjų sukuriamas vienas „string“ kintamasis, kur kiekvienas koordinačių poros kintamasis yra atskirtas vertikaliu brūkšniu „|“ (33 eilutė).
7. Vykdomas ciklas, kuris kiekvienai žingsnio koordinačių kiekio „string“ kintamajam sukuria URL užklausą „Jawg.io“ API, kuri susideda iš koordinačių kintamojo, pagrindinio URL ir vartotojo prieigos žymės (angl. *access token*) (35–40 ir 93 eilutės).
8. Kiekvienam žingsniui yra siunčiama HTTP užklausa, kuri grąžina JSON duomenų struktūrą su koordinatėmis ir jų aukščiais virš jūros lygio (41–44 eilutės).
9. Iš JSON duomenų struktūros atskiriamos aukščio virš jūros lygio reikšmės ir jos yra patalpinamos į masyvą atitinkamu indeksu (45 eilutė).
10. Kas 100 žingsnių yra atspausdinamas atliktų žingsnių kiekis (naudotas žingsnio dydis – 100 įrašų, todėl spausdinama, kas 10000 įrašų), tai atliekama, kad būtų galima stebėti procesą, ypač, kadangi tinkle atliekamos operacijos yra lėtesnės ir yra ribotas „Jawg.io“ API duodamų operacijų skaičius (37 ir 47–49 eilutės).
11. Vykdomas ciklas, kuris kiekvienam nuskaitytam įrašui sukuria naują įrašą su suapvalintomis iki pasirinkto skaičiaus po kablelio reikšmėmis bei aukščio virš jūros lygio ir suapvalinta iki dešimčių aukščio virš jūros lygio reikšmėmis. Šie įrašai sujungiami su atitinkamais įrašais ir sukuriamas naujas sąrašas su šiais įrašais (51–65 ir 94 eilutės).
12. Sukuriami nauji laukų antraštės pavadinimai bei pertvarkomos naujai pridėtų laukų (suapvalintų ir aukščio reikšmių) pozicijos (47–87 ir 95 eilutės).
13. Sukuriamas naujas CSV failas nurodytu keliu ir pavadinimu (66–72 ir 96 eilutės).
14. Įrašoma nauja antraštė su papildytais pavadinimais (69 eilutė).
15. Vykdomas ciklas, kuris kiekvieną naujo įrašų sąrašo elementą įveda į jam skirtą naujo CSV failo poziciją (70–72 eilutės).
16. Išvedama operacijos užbaigimo žinutė (97 eilutė).

3.1.1.4. Papildytų duomenų paruošimas: metų ir valandų išskyrimas

Ruošiant IPSS duomenis eksperimentui jau buvo panaudota didžioji dalis originaliai paruoštų duomenų laukų (papunkčiai 3.1.1.2 ir 3.1.1.3). Tačiau dėl didesnio potencialių aiškinamųjų laukų kiekio buvo sukurtas skriptas išskirti valandas iš laiko lauko (pvz. iš „09:15:39“ gaunama

„09“), kadangi laiko laukas buvo nenaudojamas ir metus iš datos lauko (pvz. iš „2015-09-13“ gaunama „2015“). „Python“ skriptas pavadintas „AddYearAndHour.py“. Naudojantis skriptu, galima nurodyti pradinių CSV failų kelius ir naujų CSV failų kelius. Skriptas atlieka šiuos veiksmus 3 kartus kiekvienai ryšio technologijai (3G, LTE, WiMAX) (Priedas nr. 4):

1. Nuskaito CSV įrašų failą pateiktą nurodytu keliu bei atskiria antraštę nuo įrašų (6–10 ir 39 eilutės).
2. Iš antraštės surandami datos ir laiko laukų indeksai (11–14 ir 40 eilutės).
3. Vykdomas ciklas, kuris kiekvienam nuskaitytam įrašui sukuria naują įrašą į kurį yra įterpiamos naujos metų ir valandų reikšmės į joms skirtas vietas. Šie įrašai sujungiami ir sukuriamas naujas sąrašas (15–24 ir 41 eilutės).
4. Sukuriami nauji laukų antraštės pavadinimai su įterptomis naujomis metų ir valandų reikšmėmis, joms skirtose pozicijose (25–31 ir 42 eilutės).
5. Sukuriamas naujas CSV failas nurodytu keliu ir pavadinimu (32–37 ir 43 eilutės).
6. Įrašoma nauja antraštė su papildytais pavadinimais (34 eilutė).
7. Vykdomas ciklas, kuris kiekvieną naujo įrašų sąrašo elementą įveda į naują CSV failą (36–37 eilutės).
8. Išvedama operacijos užbaigimo žinutė (44 eilutė).

3.1.1.5. Papildytų duomenų paruošimas: apvalinimo variacijos

Eksperimentui atlikti jau buvo sukurta pakankamai laukų (papunkčiai 3.1.1.2, 3.1.1.3 ir 3.1.1.4), tačiau patikrinti, kuri apvalinimo aukščio ir koordinacių kombinacija yra geriausia ir paaiškina daugiausiai anomalijų, reikia sugeneruoti daugiau duomenų rinkinių. Tam buvo parašytas dar vienas „Python“ skriptas – „GenerateMoreRoundedOptions.py“. Šis skriptas leidžia pasirinkti pradinius CSV failus, norimus testuoti intervalus koordinacių apvalinimui ir aukščio virš jūros lygio apvalinimui. Šiems intervalams yra sugeneruojamas po įrašų rinkinys, kurie yra išsaugomi pasirinktu pavadinimu ir keliu. Koordinacių apvalinimas vyksta skaitmenų po kablelio kiekio keitimu $[0, inf)$, aukščio virš jūros lygio apvalinimas vyksta iki dešimčių dalių $(0, inf)$ (pvz. dalinama į 1 dalį ir apvalinama iki – „10 20 30...“; į 2 dalis – „5 10 15...“; 3 – „3,33 6,67 10...“ ir t.t.). Taigi, šis skriptas pagal duotą pradinį duomenų rinkinį generuoja kiekvienam intervalų atvejui po naują duomenų rinkinį. Skriptas atlieka šiuos veiksmus 3 kartus kiekvienai ryšio technologijai (3G, LTE, WiMAX) (Priedas nr. 6):

1. Nuskaito CSV įrašų failą pateiktą nurodytu keliu bei atskiria antraštę nuo įrašų (6–10 ir 34 eilutės).
2. Iš antraštės surandami ilgumos, platumos, aukščio virš jūros lygio bei jų suapvalintų reikš-

mių laukų indeksai (11–18 ir 35 eilutės).

3. Vykdomas ciklas cikle, kuris kiekvienai apvalinimo iki skaitmenų po kablelio ir kiekvienai aukščio virš jūros lygio apvalinimo parinkčiai sugeneruoja po atskirą failą (36–43 eilutės).
4. Kiekvienam failui sukuriamas naujas unikalus pavadinimas (25–26 ir 36 eilutės).
5. Vykdomas ciklas, kuris platumo, ilgumos ir aukščio reikšmės suapvalinta iki nurodytos reikšmės ir pakeičia senas suapvalintas reikšmes naujomis. (19–24 ir 39 eilutės).
6. Sukuriamas naujas CSV failas nurodytu keliu ir pavadinimu (27–32 ir 40 eilutės).
7. Įrašoma nauja antraštė su papildytais pavadinimais (30 eilutė).
8. Vykdomas ciklas, kuris kiekvieną atnaujinto įrašų sąrašo elementą įveda į naują CSV failą (31–32 eilutės).
9. Išvedama vienos operacijos užbaigimo žinutė (40 eilutė).
10. Išvedama ciklo užbaigimo žinutė (41 eilutė).
11. Išvedama visų operacijų užbaigimo žinutė (42 eilutė).

Eksperimentui atlikti buvo nuskaityti ir apdoroti apie 300000 įrašų, vidutiniškai 100000 per failą. TODO: pridėti citatą

3.2. Eksperimento aplinkos paruošimas

Eksperimentas buvo atliekamas naudojant „Ubuntu (64-bit)“ operacinę sistemą įdiegtą virtualioje mašinoje „Oracle VM VirtualBox“. „MacroBase SQL“ įdiegtas ir paruoštas darbui naudojantis „MacroBase“ dokumentacija.

3.2.1. Virtuali mašina eksperimentui

„MacroBase“ pateikiami pavyzdžiai ir konfigūraciniai nurodymai yra pateikiami „Linux“ operacinėms sistemoms. Dėl šios priežasties, eksperimentui atlikti buvo pasirinkta atviro kodo nemokama operacinė sistema „Ubuntu“. Dėl paprastumo buvo nuspręsta naudoti virtualią mašiną operacinei sistemai įdiegti. Atviro kodo nemokama virtuali mašina „Oracle VM VirtualBox“ buvo pasirinkta šiai užduočiai. Galutinės eksperimentui naudotos sisteminės specifikacijos:

- Realioji mašina, kurioje atliekamas eksperimentas – nešiojamasis kompiuteris „Lenovo Y50-70“.
 - Realios mašinos Operacinė sistema „**Microsoft Windows 10 Pro**“, versija 10.0.17763.
 - Realios mašinos procesorius: **i7-4720HQ**
 - Realios mašinos operatyvioji atmintis: **8GB**

- **„Oracle VM VirtualBOX“** virtuali mašina, versija **6.0.8** r130520 (Qt5.6.2).
 - Virtualios mašinos operacinės sistemos „Ubuntu“ 64 bitų versiją **„Ubuntu (64-bit)“**, versija **18.04 LTS**.
 - Virtualiai mašinai suteiktas virtualus standusis diskas: **40GB**.
 - Virtualiai mašinai skirta operatyvioji atmintis: **4GB**.

3.2.1.1. Virtualios mašinos paruošimas

Virtualios mašinos paruošimas eksperimentui:

1. Iš „Oracle VM VirtualBox“ internetinės svetainės atsisiųstas naujausias „Windows 10“ operacinei sistemai skirtas diegimo failas (<https://www.virtualbox.org/wiki/Downloads>).
2. Sekant sąrankos vedlio nurodymus įdiegta „Oracle VM VirtualBox“ virtuali mašina.
3. Iš „Ubuntu“ internetinės svetainės atsisiųstas naujausias „Ubuntu (64-bit)“ operacinės sistemos ISO failas (<https://www.ubuntu.com/download/desktop>).
4. Atidarius „Oracle VM VirtualBox“ programinę įrangą pradėtas naujos virtualios mašinos pridėjimo procesas spaudžiant mygtuką su tekstu „New“.
5. Sekant naujos virtualios mašinos pridėjimo langus pasirinktos šios parinktys:
 - (a) Operacinės sistemos tipas: „Linux“.
 - (b) Versija: „Ubuntu (64-bit)“.
 - (c) Operatyviosios atminties kiekis megabaitais: 4096MB.
 - (d) Virtualus standusis diskas: 40GB.
6. „Oracle VM VirtualBox“ pagrindiniame lange pasirinkus naujai sukurtą virtualią mašiną spausťas mygtukas su tekstu „Start“.
7. Atidarytame virtualios mašinos lange pasirinktas anksčiau atsiųstas „Ubuntu (64-bit)“ ISO failas.
8. Pasirinkta „Install“ operacija ir sekant diegimo vedlį, įdiegta „Ubuntu (64-bit)“ operacinė sistema.

3.2.2. „MacroBase“ analitinio įrankio paruošimas

„MacroBase SQL“ diegiamas pagal „MacroBase“ dokumentacijoje nurodytus žingsnius. Tačiau prieš pradedant diegimo procesą reikia įdiegti įrankius bei papildomus paketus, tam kad „MacroBase“ taisyklingai veiktų bei būtų galima atlikti diegimo procesą. Eksperimentui naudotos programinės įrangos specifikacijos:

- „MacroBase“ versija **v1.0**.

- Java JRE versija **1.8.0_212**.
- Java JDK versija **1.8.0_212** („MacroBase“ kol kas palaiko tik 1.8.0 JDK versiją).
- „Apache Maven“ versija **3.6.0**.

3.2.2.1. „MacroBase SQL“ diegimui reikalingi papildomi paketai

Sėkmingam „MacroBase SQL“ diegimo procesui naujai paruoštoje „švarioje“ „Ubuntu“ operacinėje sistemoje, reikalingi papildomi pagalbiniai paketai. Šiuos paketus įdiegti buvo naudojamos šios terminalo komandos:

1. Versijavimo kontrolės sistemos „Git“ diegimas:

```
sudo apt install git
```

2. Java projektų valdymo ir diegimo priemonės „Apache Maven“ diegimas:

```
sudo apt install maven
```

3. Java JRE diegimas:

```
sudo apt install openjdk-8-jre
```

4. Java JDK diegimas:

```
sudo apt install openjdk-8-jdk
```

3.2.2.2. „MacroBase SQL“ analitinio įrankio diegimas

Sekant „MacroBase“ dokumentacijoje nurodytus žingsnius įdiegtas „MacroBase SQL“ įrankis:

1. Atidarytas „Ubuntu“ terminalas.
2. Klonuota projekto repozitorija:

```
git clone https://github.com/stanford-futuredata/macrobase.git
```

3. Pereita į „MacroBase“ naujai sukurtą direktoriją:

```
cd macrobase
```

4. Paleistas „MacroBase SQL“ kompiliavimo ir paruošimo BASH skriptas:

```
./build.sh sql
```

5. Patikrinta ar diegimas atliktas teisingai ir „MacroBase SQL“ įrankis veikia korektiškai pasinaudojus demonstraciniais duomenimis:

```
bin/macrobase-sql -f sql/demo.sql
```

3.3. Eksperimento vykdymo eiga

Šiame poskyryje aprašoma eksperimento vykdymo eiga naudojant IPSS pateiktus 2014–2018 metų duomenų rinkinius naudojantis „MacroBase SQL“ analitinį įrankį.

3.3.1. Konfigūracija

„MacroBase SQL“ paleidimui buvo naudojama komandinė eilutė, kuriai paduodami SQL formato failai. Šiuose failuose rašomos užklausos, kuriose nurodyti tokia informacija:

- Tiriamųjų duomenų CSV failo kelias.
- Skaitomų duomenų stulpelių pavadinimai ir tipai.
- Metrikų sąrašas.
- Išskirčių kriterijai (pvz. mažiau/daugiau nei X arba tam tikra procentilio reikšmė).
- Minimali palaikymo reikšmė.
- Minimali rizikos koeficiento reikšmė.
- Naudojamo koeficiento tipas.
- Atributų sąrašas.
- Rūšiavimo kriterijus.
- Rezultatų išvedimo CSV failų kelias.

Visų eksperimentų metu, visos SQL užklausos naudojo procentilio reikšmę išskirčių kriterijui nurodyti, taip pat naudojamo koeficiento tipas buvo visur nustatytas į rizikos koeficientą „risk_ratio(COUNT(*))“. Kiekvienos užklausos metu duomenų ir išeities rezultatų failų keliai buvo keičiami. Atributų sąrašas paliktas pilnas/tuščias „ON *“(„MacroBase SQL“ naudoja visus tinkamus atributus).

Kitus kintamuosiu buvo keičiama priklausomai nuo eksperimento tipo ar tiriamų duomenų. Tačiau buvo nustatytos numatytosios (angl. *default*) reikšmės:

- Numatytoji procentilio reikšmė – **0,99**.
- Numatytoji minimali palaikymo reikšmė – **0,001**.
- Numatytoji minimali rizikos koeficiento reikšmė – **3,0**.

3.3.2. Eksperimento eiga

Šiame skyriuje aprašyta eksperimento eiga su trimis duomenų rinkiniais – IPSS pateiktais (3G, LTE, WiMAX) ir „MacroBase SQL“ pateiktais demonstraciniais. Eksperimento metu buvo naudojamas atsitiktinai paimtas vienas iš rinkinių – 3G. Su juo buvo atliekami eksperimentai ir analizuojami rezultatai. Buvo bandoma gauti kuo geresnius rezultatus, keičiant parametrus.

Galiausiai tie patys parametrai buvo pritaikyti kitiems rinkiniams.

3.3.2.1. Eksperimento eiga su pritaikytais IPSS belaidžio interneto 3G duomenimis

Šio eksperimento metu buvo naudojami 3.1.1.2 papunktyje minimi ir paruošti duomenys. Iš pritaikytų IPSS 3G duomenų „MacroBase SQL“ automatiškai atrinko atributus: datą, technologiją ir operatorių. Eksperimento metu nustatyta metrika – sparta kb/s. Pasitelkus pateiktus atributus ir metrikas galima rasti kokie interneto operatoriai ar kokia technologija yra mažiau patikima, fiksuoja daugiau neįprastų duomenų savo tinkluose.

Buvo keistos procentilio, minimalios rizikos koeficiento ir minimalios palaikymo reikšmės. Numatytosios reikšmės paminėtos punkte 3.3.1.

Eksperimento atvejai, minimi tik pakitimai nuo numatytosios reikšmės:

1. Su numatytomis reikšmėmis.
2. Procentilio reikšmė – 0,95.
3. Palaikymo reikšmė – 0,01.
4. Procentilio reikšmė – 0,95; palaikymo reikšmė – 0,01.
5. Rizikos koeficientas – 10,0; palaikymo reikšmė – 0,01.
6. Rizikos koeficientas – 50,0; palaikymo reikšmė – 0,01.

Kiekviena užklausa rezultatus įkėlė į jai priskirtą išeities rezultatų failą CSV formatu. Šiame CSV faile, laukai susidėjo iš atributų naudotų metrikų aiškinimui (šiuo atveju: datos, operatoriaus ir technologijos), palaikymo, rizikos koeficiento, išskirčių skaičiaus ir išskirtinių atributų kombinacijų sąrašo dydžio.

3.3.2.2. Eksperimento eiga su pritaikytais IPSS belaidžio interneto 3G duomenimis, papildytais aukščio virš jūros lygio ir suapvalintų reikšmių laukais

Šio eksperimento metu buvo naudojami 3.1.1.3 papunktyje aprašyti papildyti IPSS 3G duomenų rinkinys. Duomenys yra papildyti aukščio virš jūros lygio reikšmės lauku, kuris buvo gautas iš koordinačių. Taip pat buvo pridėti koordinačių ir aukščio suapvalintų reikšmių laukais. Suapvalintų laukų reikšmės buvo interpretuojamos kaip „string“ tipo kintamieji, todėl „MacroBase SQL“ galėjo jas naudoti aiškinimui.

Šio eksperimento metu buvo atliekami tie patys atvejai, kokie buvo minimi 3.3.2.1 papunktyje, tačiau buvo daugiau atributų, kuriuos „MacroBase SQL“ galėjo naudoti aiškinimui: data, operatorius, technologijos, suapvalinta ilgumos, suapvalinta platumos, suapvalinta aukščio virš jūros lygio. Metrika naudota ta pati – sparta kb/s. Rezultatai saugoti į atitinkamai pavadintus CSV duomenų failus.

3.3.2.3. Eksperimento eiga su pritaikytais IPSS belaidžio interneto 3G duomenimis, papildytais metų ir valandų laukais

Eksperimentas atliktas tuo pačiu principu, kaip 3.3.2.1 ir 3.1.1.3 papunkčiuose aprašyti eksperimentai. Tačiau atliktas su duomenimis paruoštais ir aprašytais 3.1.1.4 papunktyje. Turėti du papildomi atributai: valandos ir metai. Metrika – sparta. Rezultatai saugoti į CSV duomenų failus.

3.3.2.4. Eksperimento eiga su pritaikytais IPSS belaidžio interneto 3G duomenimis, tikrinami optimaliausi reikšmių apvalinimai

Eksperimentas atliktas naudojant duomenų rinkinius aprašytus 3.1.1.5 papunktyje, tačiau atributai ir metrikos bei jų skaičius nepakito nuo 3.3.2.3 papunktyje naudotų. Šio eksperimento esmė buvo iš sugeneruotų duomenų rinkinių surasti labiausiai naudingą.

Eksperimento atvejai:

1. Suapvalintų koordinatų skaitmenų skaičius po kablelio – 3; aukštis apvalinamas iki dešimčių (pvz. 7,6 suapvalinamas iki 10; 6,3 – iki 10).
2. Suapvalintų koordinatų skaitmenų skaičius po kablelio – 2; aukštis apvalinamas iki dešimčių.
3. Suapvalintų koordinatų skaitmenų skaičius po kablelio – 1; aukštis apvalinamas iki dešimčių.
4. Suapvalintų koordinatų skaitmenų skaičius po kablelio – 0; aukštis apvalinamas iki dešimčių.
5. Suapvalintų koordinatų skaitmenų skaičius po kablelio – 1; aukštis apvalinamas iki pusių – penketų (pvz. 7,6 suapvalinamas iki 10; 6,3 – iki 5).
6. Suapvalintų koordinatų skaitmenų skaičius po kablelio – 1; aukštis apvalinamas iki trečdalių (pvz. 7,6 suapvalinamas iki 6,67; 6,3 iki – 6,67).
7. Suapvalintų koordinatų skaitmenų skaičius po kablelio – 0; aukštis apvalinamas iki pusių.
8. Suapvalintų koordinatų skaitmenų skaičius po kablelio – 0; aukštis apvalinamas iki trečdalių.

- 3.3.2.5. Eksperimento eiga su pritaikytais IPSS belaidžio interneto LTE, pritaikius išanalizuotus 3G parametrus
- 3.3.2.6. Eksperimento eiga su pritaikytais IPSS belaidžio interneto WiMAX, pritaikius išanalizuotus 3G parametrus
- 3.3.3. Eksperimentinio tyrimo rezultatų analize

Rezultatai ir išvados

Rezultatų ir išvadų dalyje išdėstomi pagrindiniai darbo rezultatai (kažkas išanalizuota, kažkas sukurta, kažkas įdiegta), toliau pateikiamos išvados (daromi nagrinėtų problemų sprendimo metodų palyginimai, siūlomos rekomendacijos, akcentuojamos naujovės). Rezultatai ir išvados pateikiami sunumeruotų (gali būti hierarchiniai) sąrašų pavidalu. Darbo rezultatai turi atitikti darbo tikslą.

Literatūra

- [Ast16] Anthony Asta. Observability at twitter: technical overview, part i. 2016. URL: https://blog.twitter.com/engineering/en_us/a/2016/observability-at-twitter-technical-overview-part-i.html (tikrinta 2019-03-10).
- [BGM⁺17] Peter Bailis, Edward Gan, Samuel Madden, Deepak Narayanan, Kexin Rong ir Sahaana Suri. Macrobases: prioritizing attention in fast data. *Proceedings of the 2017 ACM International Conference on Management of Data*, SIGMOD '17, p. 541–556, Chicago, Illinois, USA. ACM, 2017. ISBN: 978-1-4503-4197-4. DOI: 10.1145/3035918.3035928. URL: <http://doi.acm.org/10.1145/3035918.3035928>.
- [BGR⁺17a] Peter Bailis, Edward Gan, Kexin Rong ir Sahaana Suri. Demonstration: macrobase, a fast data analysis engine. *Proceedings of the 2017 ACM International Conference on Management of Data*, SIGMOD '17, p. 1699–1702, Chicago, Illinois, USA. ACM, 2017. ISBN: 978-1-4503-4197-4. DOI: 10.1145/3035918.3056446. URL: <http://doi.acm.org/10.1145/3035918.3056446>.
- [BGR⁺17b] Peter Bailis, Edward Gan, Kexin Rong ir Sahaana Suri. Prioritizing attention in fast data: principles and promise. *CIDR 2017, 8th Biennial Conference on Innovative Data Systems Research*. Stanford Infolab, 2017.
- [EMC14] Dell EMC. The digital universe of opportunities: rich data and the increasing value of the internet of things. 2014. URL: <http://www.emc.com/leadership/digital-universe/>.
- [Mak17] Aybeniz S. Aliyeva Makrufa Sh. Hajirahimova. About big data measurement methodologies and indicators. *I.J.Modern Education and Computer Science*, p. 1–9, 2017-10. DOI: 10.5815/ijmecs.2017.10.01.
- [PFT⁺15] Tuomas Peltkonen, Scott Franklin, Justin Teller, Paul Cavallaro, Qi Huang, Justin Meza ir Kaushik Veeraraghavan. Gorilla: a fast, scalable, in-memory time series database. *VLDB Endowment - Proceedings of the 41st International Conference on Very Large Data Bases*, p. 1816–1827, Kohala Coast, Hawaii. Facebook Inc., 2015. (Tikrinta 2019-03-10).
- [rrtar] Lietuvos Respublikos ryšių reguliavimo tarnyba. Interneto prieigos stebėsenos sistema 2014-2018 metais. URL: http://matavimai.rrt.lt/about_archive.html (tikrinta 2019-05-20).

- [Rus18] Jokūbas Rusakevičius. *Didelių duomenų srautų analizė, anomalijų aptikimas*. Programų Sistemų Kursinis darbas, Vilniaus Universitetas, Matematikos ir Informatikos Fakultetas, Vilnius, 2018.
- [Sim71] Herbert A. Simon. Designing organizations for an information-rich world. in computers, communications, and the public interest. Martin Greenberger, redaktorius, *Computers, Communication, and the Public Interest*, p. 37–72, 1971.
- [Smo14] Sandy Smolan. The human face of big data. Trumpametražis dokumentinis filmas, 2014-04-10. URL: <https://humanfaceofbigdatafilm.com/> (tikrinta 2019-05-15).
- [Sur17] Sahaana Suri. Macrobases: a search engine for fast data streams. Presentation, DataEngConf NYC '17, 2017. URL: <http://www.dataengconf.com/macrobases-a-searchengine-%20for-fast-data-streams> (tikrinta 2019-05-21).
- [Tec19] Techopedia. Anomaly detection. 2019. URL: <https://www.techopedia.com/definition/30297/anomaly-detection> (tikrinta 2019-03-07).
- [Woo15] Alex Woodie. Kafka tops 1 trillion messages per day at linkedin. 2015. URL: <https://www.datanami.com/2015/09/02/kafka-tops-1-trillion-messages-per-day-at-linkedin/> (tikrinta 2019-03-10).

Santrumpos

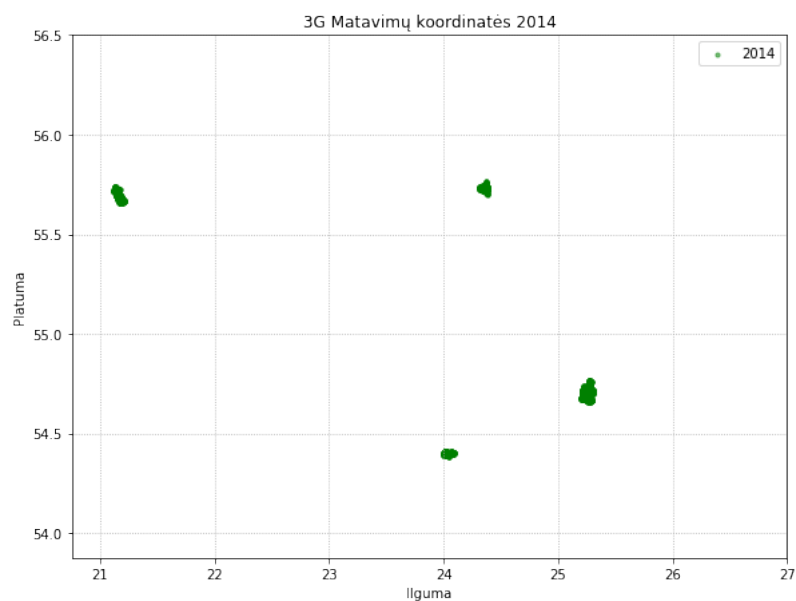
- LRTC
- CSV
- 3G
- LTE
- MAC adresas
- IPSS
- URL
- API
- HTTP
- JSON
- JRE
- JDK
- BASH
- SQL

Sąvokų apibrėžimai ir santrumpų sąrašas sudaromas tada, kai darbo tekste vartojami specialūs paaiškinimo reikalaujantys terminai ir rečiau sutinkamos santrumpos.

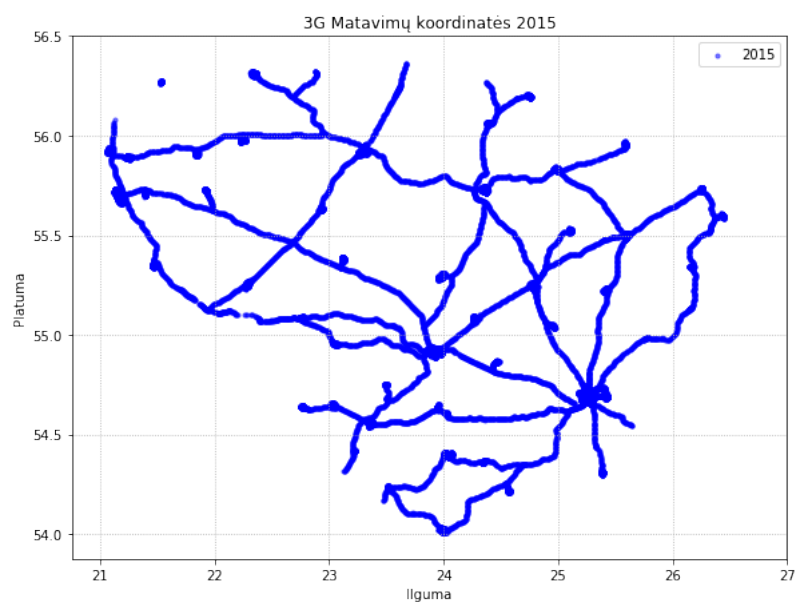
Priedas nr. 1

IPSS matavimų koordinatės 2014–2018 metais

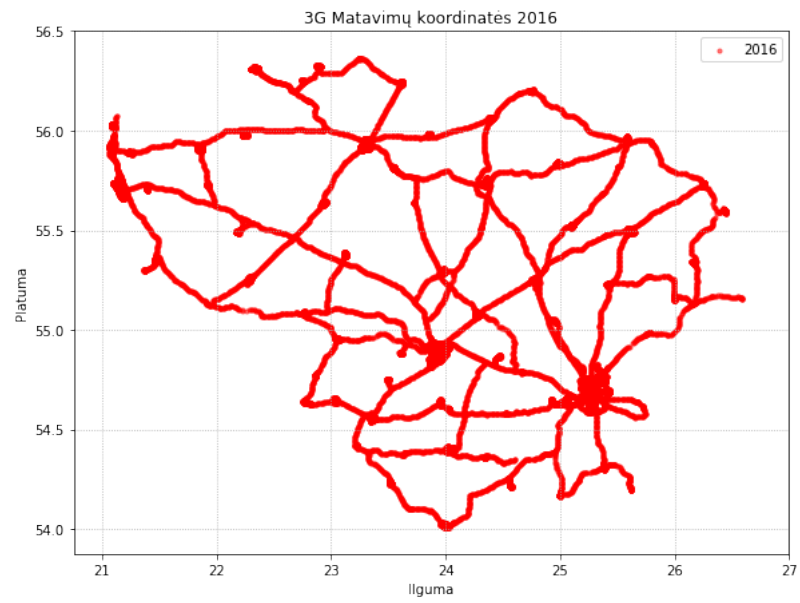
IPSS 3G ryšio technologijos matavimai 2014–2018 metais. Viso 136003 įrašai.



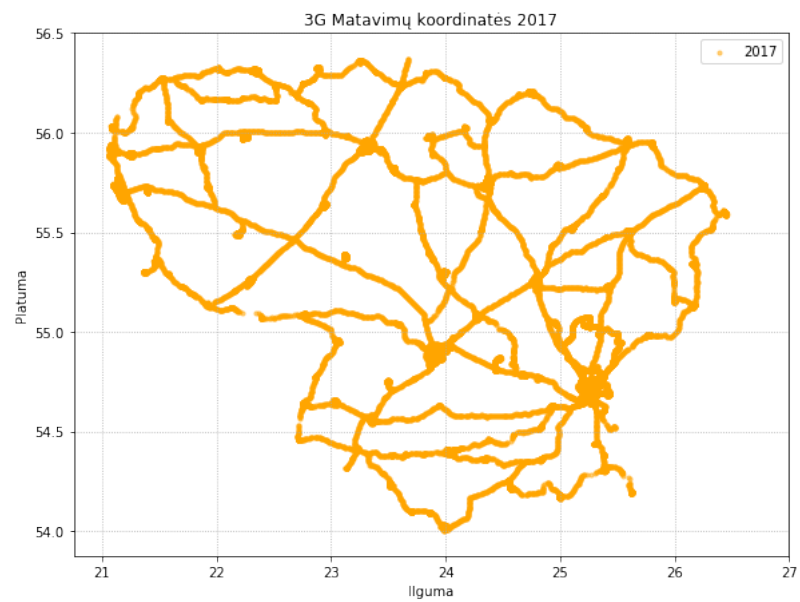
4 pav. 3G Matavimai 2014 metais. Viso – 2245 matavimai



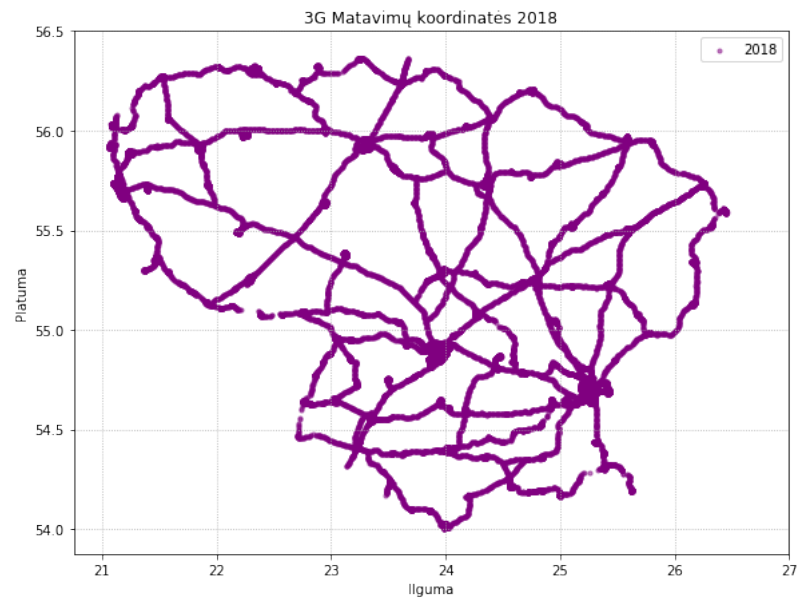
5 pav. 3G Matavimai 2015 metais. Viso – 19427 matavimai



6 pav. 3G Matavimai 2016 metais. Viso – 34769 matavimai

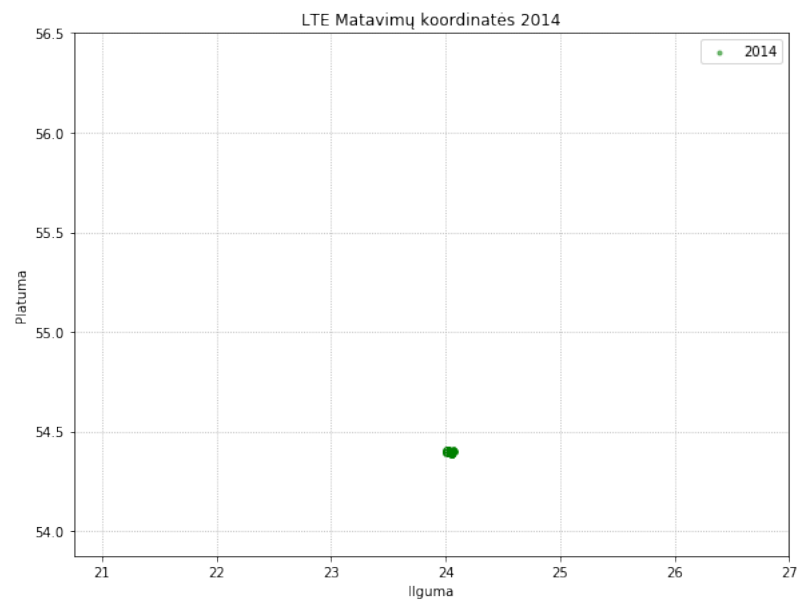


7 pav. 3G Matavimai 2017 metais. Viso – 38789 matavimai

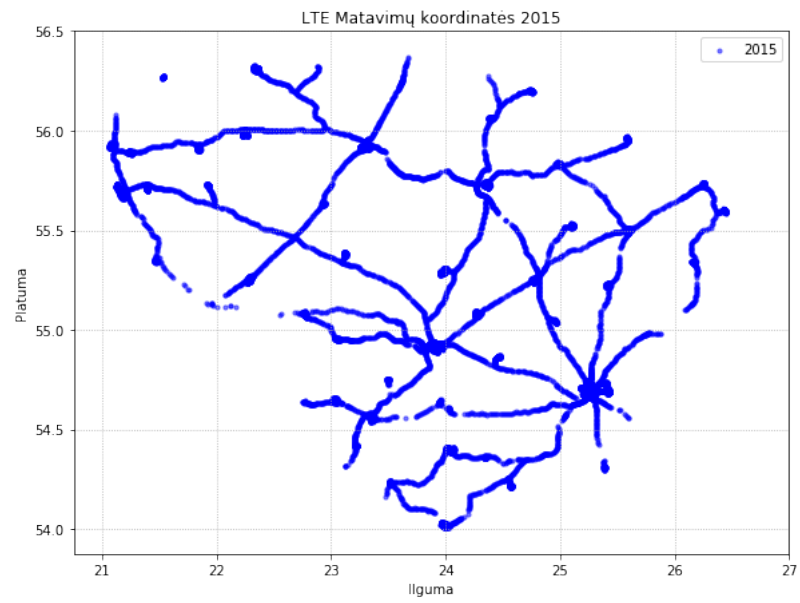


8 pav. 3G Matavimai 2018 metais. Viso – 40773 matavimai

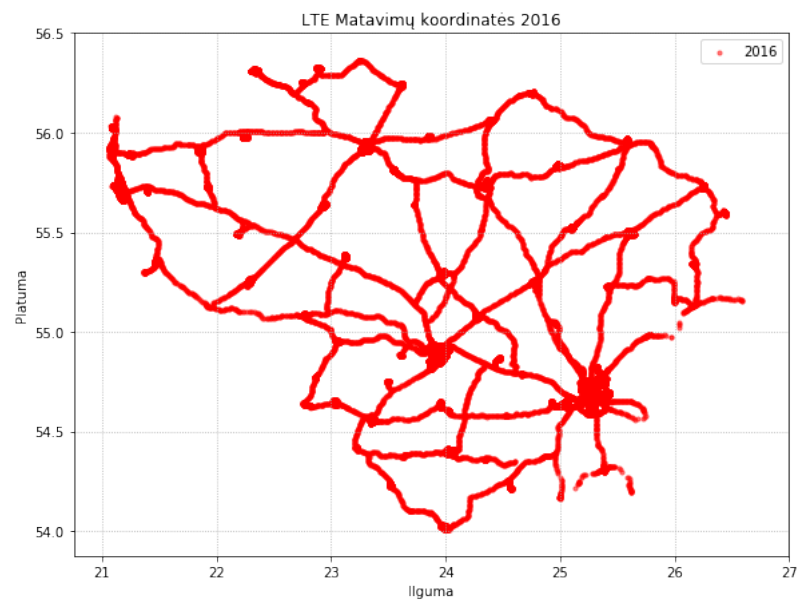
IPSS **LTE** ryšio technologijos matavimai 2014–2018 metais. Viso 157522 įrašai.



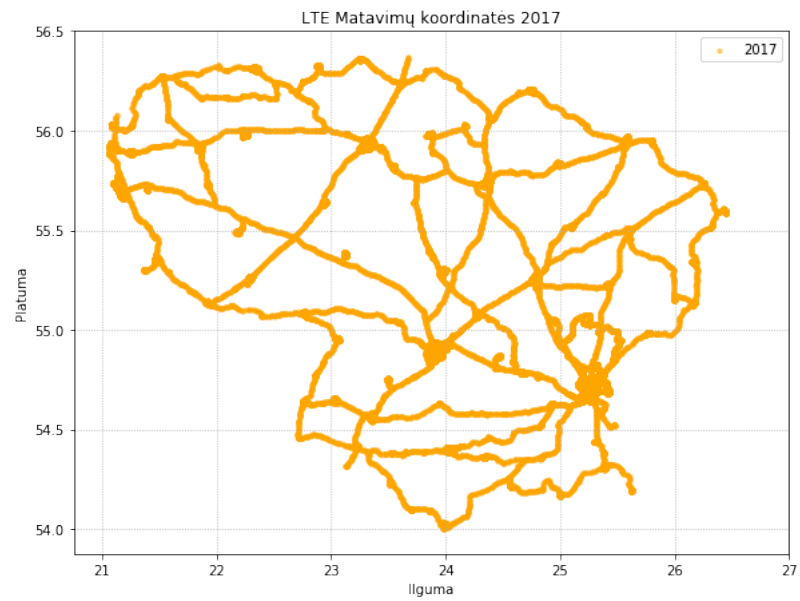
9 pav. LTE Matavimai 2014 metais. Viso – 232 matavimai



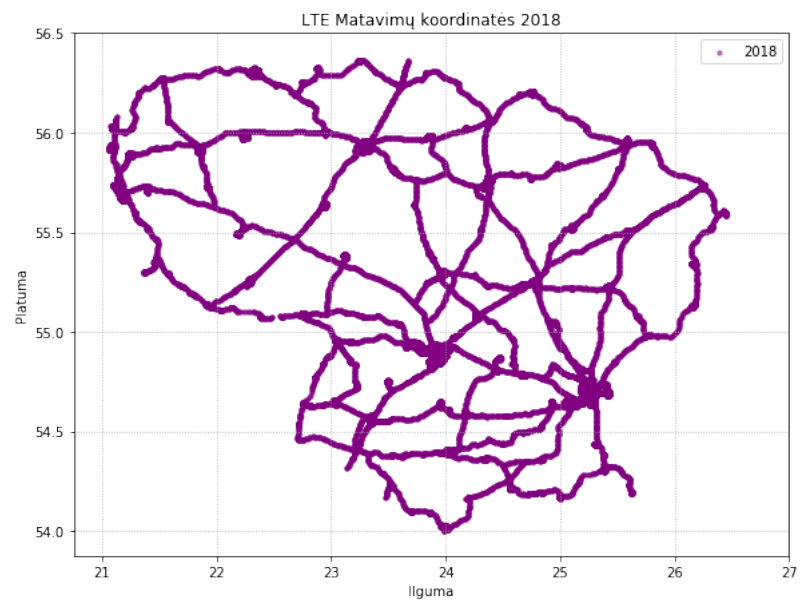
10 pav. LTE Matavimai 2015 metais. Viso – 13962 matavimai



11 pav. LTE Matavimai 2016 metais. Viso – 40738 matavimai

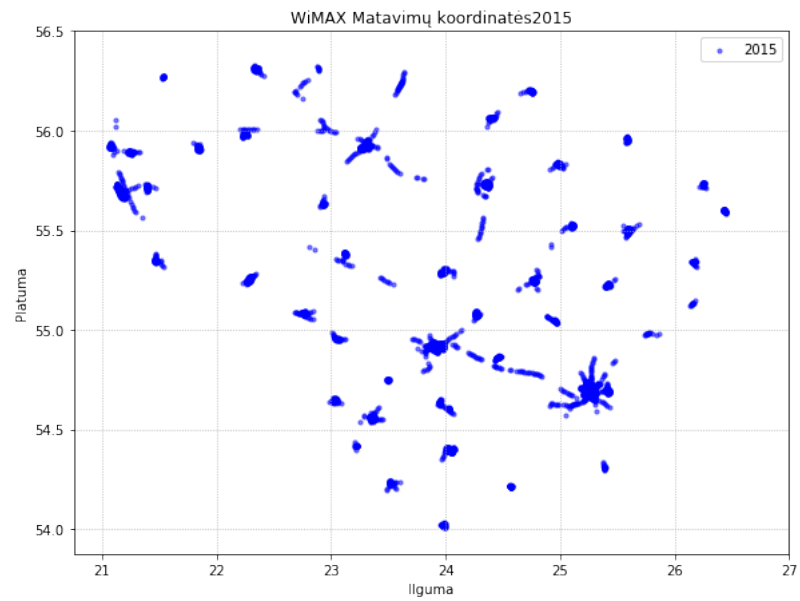


12 pav. LTE Matavimai 2017 metais. Viso – 49672 matavimai

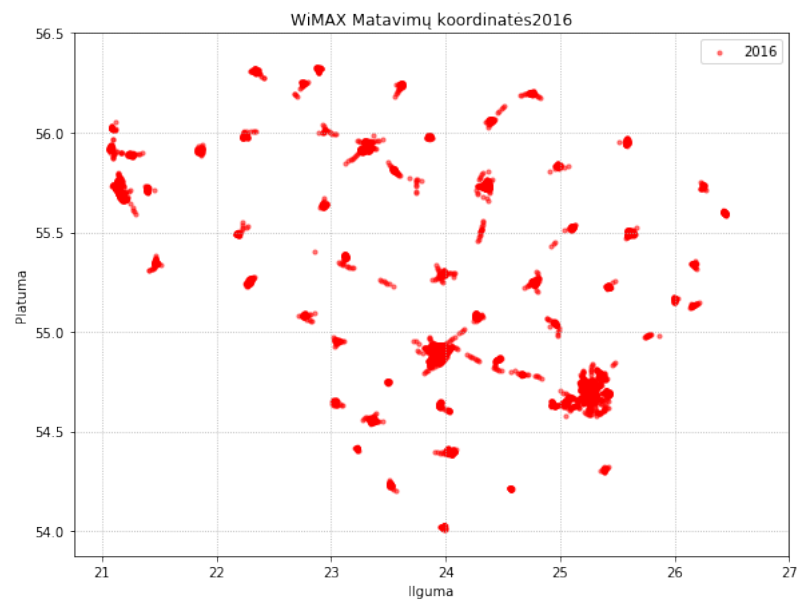


13 pav. LTE Matavimai 2018 metais. Viso – 52918 matavimai

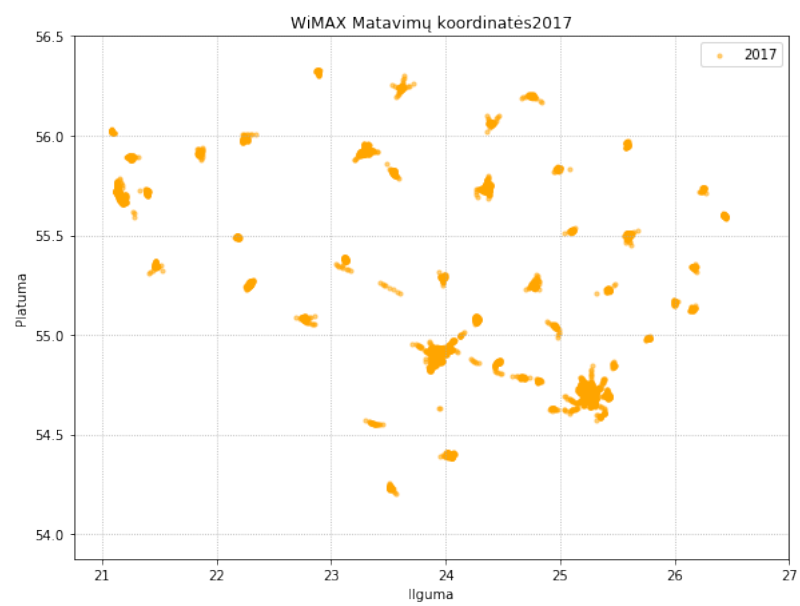
IPSS **WiMAX** ryšio technologijos matavimai 2014–2018 metais. Viso 18166 įrašai.



14 pav. WiMAX Matavimai 2015 metais. Viso – 4307 matavimai



15 pav. WiMAX Matavimai 2016 metais. Viso – 7653 matavimai



16 pav. WiMAX Matavimai 2017 metais. Viso – 6206 matavimai

Priedas nr. 2

Pradinių IPSS CSV failų apdorojimo „Python“ skriptas

„Python“ skriptas „RenameAndRemoveFields.py“ skirtas iš duotųjų IPSS duomenų CSV failų atrinkti nereikalingus laukus ir pervadinti lietuviškus antraštės pavadinimus į vieno žodžio (vienos simbolių eilutės) angliškus pavadinimus bei sugeneruoti naujus CSV failus:

```
1 #!/usr/bin/env python
2 # coding: utf-8
3 import csv
4 import numpy as np
5 import pandas as pd
6 file_3G_csv = 'Data/Raw/Matavimai-3G.csv'
7 file_LTE_csv = 'Data/Raw/Matavimai-LTE.csv'
8 file_WiMAX_csv = 'Data/Raw/Matavimai-WiMAX.csv'
9 new_file_3G_csv = 'Data/Raw/raw-3G.csv'
10 new_file_LTE_csv = 'Data/Raw/raw-LTE.csv'
11 new_file_WiMAX_csv = 'Data/Raw/raw-WiMAX.csv'
12 new_3G_labels = ['date', 'time', 'latitude', 'longitude', 'provider', 'cell_id',
13                  'rssi', 'technology', 'speed']
14 new_LTE_labels = ['date', 'time', 'latitude', 'longitude', 'provider', 'cell_id', 'rssi', 'speed']
15 new_WiMAX_labels = ['date', 'time', 'latitude', 'longitude', 'base_station_id', 'rssi', 'speed']
16 useful_3G_cols = [0, 1, 2, 3, 4, 5, 7, 8]
17 useful_LTE_cols = [0, 1, 2, 3, 4, 5, 10]
18 useful_WiMAX_cols = [0, 1, 2, 4, 5, 8]
19 def read_csv(file_name):
20     read = pd.read_csv(file_name)
21     x = np.array(read)
22     return x
23 def create_new_file(file_name, x, useful_cols, new_labels):
24     line_count = 0
25     with open(file_name, mode='w') as csv_file:
26         writer = csv.writer(csv_file, delimiter=',', quotechar='"', quoting=csv.QUOTE_MINIMAL)
27         writer.writerow(new_labels)
28         line_count = 0
29         for row in x:
30             useful_row = [row[i] for i in useful_cols]
```

```

31         useful_row[0] = date
32         useful_row.insert(1, time)
33         writer.writerow(useful_row)
34         line_count += 1
35     print(f'Processed {line_count} lines.')
36     return line_count
37 def read_and_write(file_name, new_file_name, new_labels, useful_cols):
38     x = read_csv(file_name)
39     create_new_file(new_file_name, x, useful_cols, new_labels)
40     print('Done')
41 read_and_write(file_3G_csv, new_file_3G_csv, new_3G_labels, useful_3G_cols)
42 read_and_write(file_LTE_csv, new_file_LTE_csv, new_LTE_labels, useful_LTE_cols
43 )
44 read_and_write(file_WiMAX_csv, new_file_WiMAX_csv, new_WiMAX_labels,
45 useful_WiMAX_cols)

```

Priedas nr. 3

Atrinktų IPSS CSV failų duomenų papildymo aukščio virš jūros lygio ir suapvalintų reikšmių laukais „Python“ skriptas

„Python“ skriptas skirtas jau atrinktiems CSV failams pridėti papildomus laukus: aukščio virš jūros lygio (gaunamas iš koordinatų), suapvalintų koordinatų reikšmių bei suapvalintos aukščio reikšmės. Šiems duomenims yra sugeneruojamas nauji CSV failai:

```
1  #!/usr/bin/env python
2  # coding: utf-8
3  import csv
4  import numpy as np
5  import pandas as pd
6  import urllib.parse as up
7  import urllib.request as ur
8  import json
9  def read_csv(file_name):
10     read = pd.read_csv(file_name)
11     X = np.array(read)
12     labels = np.array(read.columns)
13     return X, labels
14 def get_coord_indexes(columns):
15     lat = np.where(columns == 'latitude')[0][0]
16     lng = np.where(columns == 'longitude')[0][0]
17     return lat, lng
18 def get_coords(x, lat_index, lng_index):
19     lats = x[:, lat_index]
20     lngs = x[:, lng_index]
21     return lats, lngs
22 def create_formated_location_steps(lats, lngs, r_step):
23     prev_r = 0
24     l = len(lats)
25     locations = []
26     for r in range(r_step, l + r_step, r_step):
27         locs = []
28         for i in range(prev_r, r):
29             if i >= l:
30                 break
31             locs.append(str(lats[i]) + ',' + str(lngs[i]))
32     prev_r = r
```

```

33     locations.append("|".join(locs))
34     return locations
35 def get_elevations(locations, main_url, token):
36     elevations = np.array([])
37     a = 0
38     for loc in locations:
39         params = {'locations': loc, 'access-token': token}
40         full_url = main_url + up.urlencode(params)
41         request = ur.urlopen(full_url)
42         contents = request.read()
43         request.close()
44         jdata = json.loads(contents)
45         el = np.array([j['elevation'] for j in jdata])
46         elevations = np.concatenate((elevations, el))
47         a += 1
48         if a % 100 == 0:
49             print(f'Processed altitude requests: {a}')
50     return elevations
51 def create_new_X(x, lat, lng, els, c_round_digits, alt_round_parts=1):
52     new_X = []
53     for i in range(len(x)):
54         item = x[i]
55         new_item = []
56         new_item.append(round(item[lat], c_round_digits))
57         new_item.append(round(item[lng], c_round_digits))
58         el = els[i]
59         new_item.append(el)
60         round_el = round(round(el / 10 * alt_round_parts) * 10 /
alt_round_parts, 2)
61         new_item.append(round_el)
62         complete_item = np.append(np.array(item), np.array(new_item))
63         new_X.append(complete_item)
64     new_X = np.array(new_X)
65     return new_X
66 def write_csv(file_name, x, labels, positinos):
67     with open(file_name, mode='w') as new_csv_file:
68         writer = csv.writer(new_csv_file, delimiter=',', quotechar='"',
quoting=csv.QUOTE_MINIMAL)
69         writer.writerow(labels)
70         for item in x:

```



```

71         correct_pos_item = [item[p] for p in positions]
72         writer.writerow(correct_pos_item)
73 def create_new_labels_and_positions(labels, lat, lng):
74     positions = np.arange(len(labels)).tolist()
75     labels_list = labels.tolist()
76     last_i = len(labels)
77     i = max(lat, lng)+1
78     positions.insert(i, last_i+2)
79     positions.insert(i+1, last_i)
80     positions.insert(i+2, last_i+1)
81     positions.insert(i+3, last_i+3)
82     labels_list.insert(i, 'altitude')
83     labels_list.insert(i+1, 'lat')
84     labels_list.insert(i+2, 'lng')
85     labels_list.insert(i+3, 'alt')
86     labels_list = np.array(labels_list)
87     return positions, labels_list
88 def create_new_csv_file(file_name, new_file_name, request_step, token,
89     main_url, round_digits=2, round_altitude_parts=1):
89     X, labels = read_csv(file_name)
90     lat, lng = get_coord_indexes(labels)
91     lats, lngs = get_coords(X, lat, lng)
92     locations = create_formated_location_steps(lats, lngs, request_step)
93     els = get_elevations(locations, main_url, token)
94     new_X = create_new_X(X, lat, lng, els, round_digits, round_altitude_parts)
95     positions, new_labels = create_new_labels_and_positions(labels, lat, lng)
96     write_csv(new_file_name, new_X, new_labels, positions)
97     print('Done')
98 file_3G_csv = 'Data/Raw/raw-3G.csv'
99 file_LTE_csv = 'Data/Raw/raw-LTE.csv'
100 file_WiMAX_csv = 'Data/Raw/raw-WiMAX.csv'
101 new_file_3G_csv = 'Data/Processed/altitude-11-3G.csv'
102 new_file_LTE_csv = 'Data/Processed/altitude-11-LTE.csv'
103 new_file_WiMAX_csv = 'Data/Processed/altitude-11-WiMAX.csv'
104 token = "nIXn7WWJTSeguGlaEuSioR2hM5PqAZUcRfGqCOV9RJvDwVJIYwDdDmwVDKvs5FYm"
105 main_url = "https://api.jawg.io/elevations?"
106 request_step = 100
107 round_digits = 1
108 round_altitude = 1
109 create_new_csv_file(file_3G_csv, new_file_3G_csv, request_step, token,

```

```
        main_url, round_digits, round_altitude)
110 create_new_csv_file(file_LTE_csv, new_file_LTE_csv, request_step, token,
        main_url, round_digits, round_altitude)
111 create_new_csv_file(file_WiMAX_csv, new_file_WiMAX_csv, request_step, token,
        main_url, round_digits, round_altitude)
```

Priedas nr. 4

Atrinktų IPSS CSV failų duomenų papildymo metų ir valandų laukais „Python“ skriptas

„Python“ skriptas skirtas jau atrinktiems CSV failams pridėti papildomus laukus: metų ir valandų (gaunamas iš datos ir laiko). Šiems duomenims yra sugeneruojamas nauji CSV failai:

```
1 #!/usr/bin/env python
2 # coding: utf-8
3 import csv
4 import numpy as np
5 import pandas as pd
6 def read_csv(file_name):
7     read = pd.read_csv(file_name)
8     X = np.array(read)
9     labels = np.array(read.columns)
10    return X, labels
11 def get_datetime_indexes(columns):
12     dd = np.where(columns == 'date')[0][0]
13     tt = np.where(columns == 'time')[0][0]
14     return dd, tt
15 def create_new_x(x, dd, tt):
16     i = max(dd, tt)
17     new_x = []
18     for row in x:
19         new_row = row.tolist()
20         new_row.insert(i+1, row[dd][:4])
21         new_row.insert(i+2, row[tt][:2])
22         new_x.append(new_row)
23     new_x = np.array(new_x)
24     return new_x
25 def create_new_l(l, dd, tt, yl, hl):
26     i = max(dd, tt)
27     new_l = l.tolist()
28     new_l.insert(i+1, yl)
29     new_l.insert(i+2, hl)
30     new_l = np.array(new_l)
31     return new_l
32 def write_csv(file_name, x, labels):
33     with open(file_name, mode='w') as new_csv_file:
```

```

34     writer = csv.writer(new_csv_file , delimiter=',', quotechar='"',
    quoting=csv.QUOTE_MINIMAL)
35     writer.writerow(labels)
36     for row in x:
37         writer.writerow(row)
38 def your_hour_change_file(file_name , new_file_name , year_label , hour_label):
39     x, l = read_csv(file_name)
40     dd, tt = get_datetime_indexes(l)
41     new_x = create_new_x(x, dd, tt)
42     new_l = create_new_l(l, dd, tt , year_label , hour_label)
43     write_csv(new_file_name , new_x, new_l)
44     print('Done')
45 file_3G_csv = 'Data/Processed/altitude-3G.csv'
46 file_LTE_csv = 'Data/Processed/altitude-LTE.csv'
47 file_WiMAX_csv = 'Data/Processed/altitude-WiMAX.csv'
48 new_file_3G_csv = 'Data/Processed/yh-3G.csv'
49 new_file_LTE_csv = 'Data/Processed/yh-LTE.csv'
50 new_file_WiMAX_csv = 'Data/Processed/yh-WiMAX.csv'
51 year_label = 'year'
52 hour_label = 'hour'
53 your_hour_change_file(file_3G_csv , new_file_3G_csv , year_label , hour_label)
54 your_hour_change_file(file_LTE_csv , new_file_LTE_csv , year_label , hour_label)
55 your_hour_change_file(file_WiMAX_csv , new_file_WiMAX_csv , year_label ,
    hour_label)

```

Priedas nr. 5

Atrinktų IPSS papildomų CSV failų generavimo „Python“ skriptas

„Python“ skriptas skirtas jau atrinktiems CSV failams sugeneruoti skirtingų koordinacių ir aukščio virš jūros lygio reikšmių apvalinimo duomenų rinkinius. Šiems duomenims yra sugeneruojamas nauji CSV failai:

```
1 #!/usr/bin/env python
2 # coding: utf-8
3 import csv
4 import numpy as np
5 import pandas as pd
6 def read_csv(file_name):
7     read = pd.read_csv(file_name)
8     X = np.array(read)
9     labels = np.array(read.columns)
10    return X, labels
11 def get_col_indexes(columns):
12     latitude = np.where(columns == 'latitude')[0][0]
13     longitude = np.where(columns == 'longitude')[0][0]
14     altitude = np.where(columns == 'altitude')[0][0]
15     lat = np.where(columns == 'lat')[0][0]
16     lng = np.where(columns == 'lng')[0][0]
17     alt = np.where(columns == 'alt')[0][0]
18     return latitude, longitude, altitude, lat, lng, alt
19 def change_x(x, latitude, longitude, altitude, lat, lng, alt, rd, ra):
20     for i in range(len(x)):
21         item = x[i]
22         item[lat] = round(item[latitude], rd)
23         item[lng] = round(item[longitude], rd)
24         item[alt] = round(round(item[altitude] / 10 * ra) * 10 / ra, 2)
25 def name_new_file(file_name, rd, ra):
26     return file_name.format(rd, ra)
27 def write_csv(file_name, x, labels):
28     with open(file_name, mode='w') as new_csv_file:
29         writer = csv.writer(new_csv_file, delimiter=',', quotechar='\"',
30                               quoting=csv.QUOTE_MINIMAL)
31         writer.writerow(labels)
32         for row in x:
33             writer.writerow(row)
34 def perform_all_tests(file_name, new_file_name, rdt, rat):
35     x, l = read_csv(file_name)
```

```

35     latitude , longitude , altitude , lat , lng , alt = get_col_indexes(1)
36     for rd in rdt:
37         for ra in rat:
38             full_name = name_new_file(new_file_name , rd , ra)
39             change_x(x, latitude , longitude , altitude , lat , lng , alt , rd , ra)
40             write_csv(full_name , x , 1)
41             print(f'{rd}{ra} Done.')
42         print('Cycle.')
43     print('All done.')
44 file_3G_csv = 'Data/Processed/yh-3G.csv'
45 file_LTE_csv = 'Data/Processed/yh-LTE.csv'
46 file_WiMAX_csv = 'Data/Processed/yh-WiMAX.csv'
47 new_file_3G_csv = 'Data/Processed/3G/{}/{}-3G.csv'
48 new_file_LTE_csv = 'Data/Processed/LTE/{}/{}-LTE.csv'
49 new_file_WiMAX_csv = 'Data/Processed/WiMAX/{}/{}-WiMAX.csv'
50 round_digit_tests = range(3+1)
51 round_altitude_tests = range(1, 5+1)
52 perform_all_tests(file_3G_csv , new_file_3G_csv , round_digit_tests ,
53                   round_altitude_tests)
53 perform_all_tests(file_LTE_csv , new_file_LTE_csv , round_digit_tests ,
54                   round_altitude_tests)
54 perform_all_tests(file_WiMAX_csv , new_file_WiMAX_csv , round_digit_tests ,
55                   round_altitude_tests)

```

Priedas nr. 6

„MacroBase SQL“ SQL užklausų pavyzdys

Darbo metu naudotos „MacroBase SQL“ SQL užklausos pavyzdys su pradiniais duomenimis, papildytais su aukščio virš jūros lygio ir suapvalintomis reikšmėmis duomenimis, ir papildytais metų bei valandų laukais duomenimis:

```
1 #Pradiniai duomenys
2 IMPORT FROM CSV FILE '~/Data/Raw/raw-3G.csv' INTO r3G(date string, time string
  , latitude double, longitude double, provider string, cell_id string, rssi
  double, technology string, speed double);
3 SELECT * FROM
4   DIFF
5     (SPLIT (
6       SELECT *, percentile(speed) as pct FROM r3G)
7     WHERE pct > 0.99)
8   ON *
9   WITH MIN SUPPORT 0.01 MIN RATIO 3.0
10  COMPARE BY risk_ratio(COUNT(*))
11  ORDER BY support
12  INTO OUTFILE 'Output Data/3G/raw-3G.csv' FIELDS TERMINATED BY ',';
13 # Papildyti aukscio ir suapvalintomis reiksmemis duomenys
14 IMPORT FROM CSV FILE '~/Data/Processed/altitude-3G.csv' INTO a3G(date string,
  time string, latitude double, longitude double, altitude double, lat
  string, lng string, alt string, provider string, cell_id string, rssi
  double, technology string, speed double);
15 SELECT * FROM
16   DIFF
17     (SPLIT (
18       SELECT *, percentile(speed) as pct FROM a3G)
19     WHERE pct > 0.99)
20   ON *
21   WITH MIN SUPPORT 0.01 MIN RATIO 3.0
22  COMPARE BY risk_ratio(COUNT(*))
23  ORDER BY support
24  INTO OUTFILE 'Output Data/3G/altitude-3G.csv' FIELDS TERMINATED BY ',';
25 # Papildyti metu ir valandu laukais duomenys
26 IMPORT FROM CSV FILE '~/Data/Processed/yh-3G.csv' INTO y3G(date string, time
  string, year string, hour string, latitude double, longitude double,
  altitude double, lat string, lng string, alt string, provider string,
  cell_id string, rssi double, technology string, speed double);
```

```
27 SELECT * FROM
28     DIFF
29     (SPLIT (
30         SELECT *, percentile(speed) as pct FROM y3G)
31         WHERE pct > 0.99)
32 ON *
33 WITH MIN SUPPORT 0.01 MIN RATIO 3.0
34 COMPARE BY risk_ratio(COUNT(*))
35 ORDER BY support
36 INTO OUTFILE 'Output Data/3G/yh-3G.csv' FIELDS TERMINATED BY ',';
```