

LATVIJAS UNIVERSITĀTES
DATORIKAS FAKULTĀTE

**PAPILDU FUNKCIJU IZSTRĀDE
MAŠĪNMĀCĪŠANĀS TEKSTA ANALIZATORAM**

KVALIFIKĀCIJAS DARBS

Autors: Mikus Kalniņš
Studenta apliecības numurs: mk17110
Darba vadītājs: prof. Jānis Zuters

RĪGA 2019

ANOTĀCIJA

Teksta analīzes veikšanai ir pieejamas vairākas metodes, šajā dokumentā ir apskatīts kā pielietot *Correlation Explanation*[7] (korelāciju skaidrošana, turpmāk *CorEx*[7]) teksta analīzes metodi, implementējot atvērtā pirmkoda (*open-source*) bibliotēku *corextopic*[3], jau izstrādātā sistēmā - mašīnmācīšanās teksta analizatorā (turpmāk *MMTA*). Darbā ir aprakstīta *MMTA* pamatdarbība un tā mijiedarbība ar implementēto *corextopic*[3] bibliotēku.

MMTA ir programmprodukts ar implementētām vairākām bibliotēkām, kas dod iespēju programmprodukta lietotājam izvēlēties starp vairākām datu apstrādes metodēm saistībā ar teksta analīzi. *MMTA* darbība iedalās. Informācijas ekstrakcija un tēmu modelēšana ir divi *MMTA* darbības iedalījumu piemēri. Darbā tiek apskatīta tēmu modelēšana, jo *corextopic*[3] bibliotēka sniedz bagātīgi modelētas tēmas, tā raksturojot apstrādājamo datu kopumu.

Atslēgvārdi: *CorEx*, tēmu modelēšana, mašīnmācīšanās, teksta analīze, python

ABSTRACT

Additional Feature Development for Machine Learning Text Analyzer

Text analysis has many algorithms that try to solve the problem of topic modeling. This document explains how the text analysis method called *Correlation Explanation*[7] was integrated into a *Machine Learning Text Analyzer (MaLTA)* using an open-source *Python* library *corextopic*[3]. The document describes core functions of *MaLTA* and how they interact with the integrated *corextopic*[3] library.

MaLTA is a program product with many integrated open-source libraries, that offer the user of *MaLTA* to choose from a variety of methods to use for text analysis. There are several categories of text analysis that *MaLTA* works with. *Topic Modeling* and *Information Extraction* are two examples. The document describes the process of Topic Modeling since *corextopic*[3] offers richly modeled topics in such a way describing the input data.

Key-words: CorEx, topic modeling, machine learning, text analysis, Python

SATURA RĀDĪTĀJS

Anotācija	2
Abstract	3
Satura rādītājs.....	4
Terminu skaidrojums	6
Ievads	8
1. Programmatūras prasību specifikācija	9
1.1. Ievads	9
1.1.1. Nolūks	9
1.1.2. Darbības sfēra.....	9
1.1.3. Saistība ar citiem dokumentiem	9
1.1.4. Pārskats.....	9
1.2. Vispārējais apraksts	10
1.2.1. Produkta perspektīva	10
1.2.2. Produkta funkcijas	10
1.2.3. Lietotāja raksturiezīmes	12
1.2.4. Vispārēji ierobežojumi	12
1.3. Konkrētās prasības	13
1.3.1. Vispārīgā informācija	13
1.3.2. Funkcionālās prasības.....	14
1.4. Ārējā saskarne.....	33
1.4.1. Lietotāju saskarne.....	33
1.4.2. Aparatūras saskarne.....	34

1.4.3.	Programmatūras saskarne	35
1.4.4.	Sakaru saskarne	35
1.5.	Nefunkcionālās prasības	36
1.5.1.	Veikspējās prasības	36
1.6.	Konceptuālās datu plūsmas	36
1.6.1.	Sistēmas koncepcija	36
1.6.2.	Pirmā līmeņa Datu plūsmu diagramma –	37
2.	Programmatūras projektējuma apraksts	38
2.1.	Ievads	38
2.1.1.	Nolūks	38
2.1.2.	Darbības sfēra	38
2.1.3.	Saistība ar citiem dokumentiem	38
2.1.4.	Pārskats	38
2.2.	Tehniskie risinājumi	38
2.3.	Dekompozīcijas apraksts	40
2.3.1.	Objektu dekompozīcija	40
3.	Testēšanas dokumentācija	53
3.1.	Automātiskie vienībtesti un funkcionālie testi	53
3.2.	Vienībtestēšanas rezultāti	53
3.2.1.	TestCorExExtractor	54
3.2.2.	TestCorExConfiguration	67
3.2.3.	TestFactory	68
4.	Projekta organizācija	72
5.	Darbietilpības novērtējums	74
	Nobeigums	78
	Izmantotā literatūra	79

TERMINU SKAIDROJUMS

Atslēga – programmēšanas vērtības nosaukums, ko angļiski sauc par *key*, kas ir daļa no *dictionary* (vārdnīcas).

Bibliotēka – rīks, iepriekš izveidots kods, ar izstrādei lietderīgām funkcijām.

Būla tips – patiesuma vērtība. Jā vai nē, 0 vai 1, Patiess vai nepatiess. Programmēšanā parasti atsaucas uz vērtībām “True” un “False” vai 0 un 1.

CorEx – *Correlation Explanation*[7] no aģļu valodas tulkojot: korēlāciju skaidrošana.

Celmošana – teksta apstrādes process, kas apstrādā vārdus tā, lai no tiem paliktu tikai vārda sakne un nominatīvs. Bieži vien noņemot galotnes un pārveidojot vārdus vienskaitļa formā. Piemēram, vārds “Trouble” tiek pārvērst par “Troubl”.

Datu tips – datu paveids. Piemēram, skaitlisks, simbolu virkne, simbols.

Inicializē – izveido, iesāk darbību.

Instance – objekts, objekta gadījums, esamība, objekts pats. Piemēram, monēta var būt ideja, taču monētas instance varētu būt kāda konkrēta monētā, kas tiek turēta rokā.

Iterēt – veikt kādu darbību atkārtoti, secīgi.

Klase – programmēšanas procesā izveidots datu objekts, kas var saturēt funkcijas, vērtības. Veidne.

Konsole – vadības saskarne.

Lemmatizācija – līdzīgs process kā celmošana.

MMTA – saīsinājums no “Mašīnmācīšanās teksta analizators”.

Nebūtisks vārds – vārds, kas dod maz informācijas par teksta kopējo nozīmi, vai ir ar ļoti plašu nozīmi. Piemēram, “es, tu, mēs, grib, vienreiz”. Pretējs piemērs, “dators, sāp, galvas sāpes, armija”.

Noklusējuma – iepriekš izvēlēta vai uzstādīta vērtība.

PPA – Programmatūras projektējuma apraksts.

PPS – Programmatūras prasību specifikācija.

Teksta dokuments – viens teksta vienums. Piemēram, e-pasta saturs, grāmatas virsraksts, dziesmas vārdi. Šādi teksta dokumenti veido teksta dokumentu kopu.

Tokenizēšana – teksta analīzes sfērā apzīmē procesu, kas, sadala teikumu tokenos.

Tokens – teksta analīzes sfērā ar apzīmē vārdu, vai vārda sakni.

Trenēt – mašīnmācīšanās termins, kas apzīmē procesu, kad iteratīvi modeli trenē uz kādiem datiem, lai iegūtu labākus rezultātus izmantojot trenētu mašīnmācīšanās modeli.

Vārdnīca – Datu struktūra, kas nodrošina piekļuvi datiem, programmām vai datnēm, izmantojot to identifikatorus[10].

IEVADS

Kvalifikācijas darba mērķis ir izstrādāt papildus funkcionalitātes prototipa implementāciju pasūtītāja mašīnmānīšanās teksta analizatorā, tās sniedzot *MMTA* lietotājiem vairāk iespēju izvēlēties no vairākiem tēmu modelēšanas paņēmieniem analizējot lietotāja izvēlētos teksta datus.

Izstrāde ir nepieciešama, jo katra tēmu modelēšanas metode strādā pamatojoties uz dažādiem pieņēmumiem. Dodot lietotājam izvēli starp vairākām metodēm, lietotājs *MMTA* lietošanas laikā var atrast metodi, kas pēc lietotāja domām sniedz labāko rezultātu jeb lietotājs var izvēlēties metodi, kas viņprāt vislabāk raksturo analizējamus datus.

Iepriekšminēto iemeslu dēļ pasūtītājs ir izvēlējies korelāciju ekstrakcijas(*CorEx*[7]) metodi kā nākamo implementējamo metodi savā sistēmā, un šo implementāciju var veikt izmantojot atvērtā pirmkoda risinājums *corextopic*[3] integrēšanu pasūtītāja sistēmā. *Corextopic*[3] ir atvērtā pirmkoda risinājums, kas realizē *CorEx*[7] metodes darbību

Integrēšanas veikšanai tiks izmantota *Python 3.5.5* programmēšanas valoda, objektorientētas programmēšanas pieeja un atvērtā pirmkoda risinājumi kā *corextopic*[3], *CountVectorizer*, *SnowBallStemmer*, *WordNetLemmatizer* un citus.

Sākotnējās konsultācijās ar pasūtītāju tika nolemts, ka jātiek izveidotam konceptuālam prototipam, ko lēnām izstrādes laikā papildinot tiks izveidota jau pilnīgi lietojama *CorEx*[7] implemetnācija sistēmā ar tās pamatfunkcijām, kuras ir ielasīt konfigurācijas failu, ielasīt lietotāja dotos ievaddatus un kā rezultātu izvadīt trīs failus – sakategorizētus ievaddatus divos veidos un kategoriju jeb tēmu skaidrojumu.

Izpildot iepriekšējā teikumā minētās prasības, var tālāk tikt veikta *CorEx*[7] implementācijas funkcionāla papildināšana un implementētā risinājuma darbības uzlabošana.

1. PROGRAMMATŪRAS PRASĪBU SPECIFIKĀCIJA

1.1. Ievads

1.1.1. Nolūks

Programmatūras prasību specifikācija (turpmāk PPS) ir dokuments, kas raksturo prasības, kas jāievēro *corextopic*[3] bibliotēkas implementācijas veikšanai. Implementācijas mērķis ir papildināt *MMTA* tēmu modelēšanas arsenālu ar bibliotēku, kas ir spējīga teksta dokumentu datu kopā atrast iedalījumu tēmās un atpazīt kādas tēmas piemīt katram dokumentam dotajā datu kopā.

1.1.2. Darbības sfēra

Izstrādātā sistēma jeb *CorEx*[7] *implementācija* ir lielākas sistēmas (*MMTA*) funkcionalitātes papildinājums, kas nodrošina teksta dokumentu kategorizēšanu izmantojot *corextopic*[3] mašīnmācīšanās bibliotēku.

1.1.3. Saistība ar citiem dokumentiem

Dokumenta noformēšanā ievērotas standarta LVS 68:1996 „Programmatūras prasību specifikācijas ceļvedis” prasības.

1.1.4. Pārskats

PPS sastāv no šādām nodaļām:

- 1) Ievads. Apraksta PPS nolūku, darbības sfēru, saistību ar citiem dokumentiem.
- 2) Vispārējais apraksts. Satur vispārēju aprakstu par sistēmas pamatfunkcionalitāti, tās perspektīvu, sistēmas ierobežojumiem un sistēmas lietotāja raksturiezīmes.
- 3) Konkrētās prasības. Detalizēti specificētas visas prasības, kas sistēmas izstrādātājam jāievēro, lai veiktu sistēmas projektējumu.
- 4) Ārējā saskarne. Aprakstītas gan lietotāja, gan sistēmas saskarnes
- 5) Nefunkcionālās prasības. Aprakstītas sistēmas veikspējas prasības.
- 6) Konceptuālās datu plūsmas. Attēlota sistēmas darbība izmantojot diagrammas.

1.2. Vispārējais apraksts

1.2.1. Produkta perspektīva

CorEx[7] implementācija ir veikta atbilstoši pasūtītāja prasībām, mērķiem un vajadzībām, ievērojot Latvijas Republikas un Eiropas Savienības likuma prasības. Izstrādātās sistēmas funkcionālo darbību nodrošina *CorEx*[7] implementācijas saskarne ar *MMTA*. Savukārt *MMTA* ir programmprodukts, kuras darbību un lietojamību nodrošina pasūtītājs.

Izstrādātā sistēma izmanto *MMTA* iekšējās funkcijas failu ielasīšanai un rakstīšanai, sistēmas darbības konfigurēšanai, klašu un koda struktūras pārbaudei. Izmantojot *MMTA* iekšējās funkcijas, tiek nodrošināts, ka nākotnē implementēto funkcionalitāšu darbības laikā faili tiks ielasīti un rakstīti vienādā formātā un ka tiks ievērota *MMTA* projekta programmatūras koda struktūra.

1.2.2. Produkta funkcijas

CorEx[7] implementācijas funkcionalitāte ir sistematizēta pēc klasēm un pēc to nozīmes un pamatmērķa teksta analīzes processā. Zemāk ir redzamas klases, kas nodrošina sistēmas un implementētās funkcionalitātes kopēju darbību (implementācijas klases, *MMTA* klases un bibliotēkas klases). Implementācijas klases ir klases, kas ir jāizstrādā. *MMTA* klases ir jau sistēmā izstrādātas klases, kurām pēc vajadzības ir jāveic izmaiņas, lai veidotos saskarne ar implementācijas laikā izstrādātajām klasēm. Bibliotēkas klases ir klases, kuru funkcijas tiešā veidā izmanto implementācijas laikā izstrādātajās klasēs. Pēc pasūtītāja lūguma, pielietotās *MMTA* klases tiks raksturotas vispārēji. Bibliotēkas klases tiks aprakstītas vispārēji, jo to oficiālā dokumentācija ir atrodama atsaucēs.

Implementācijas klases:

- *CorExExtractor*
- *CorExConfiguration*

MMTA klases:

- *TopicModeling*
- *Trainer*
- *FileHelper*

- *BaseExtractor*
- *Factory*
- *RunConfiguration*

Bibliotēkas bibliotēkas klase:

- *CorExModel*[3]

Tālāk seko katras klases darbības apraksts:

TopicModeling: Uzsāk *MMTA* darbību izmantojot *Trainer* klasi.

Trainer: Izmantojot klases *FileHelper* un *Factory*, ielasa konfigurācijas failu, izveido *CorExExtractor* instanci un izsauc *CorExExtractor* vajadzīgās funkcijas vajadzīgajā secībā, lai tiktu veikta teksta analīze.

FileHelper: Definē un veic dažādu formātu failu ielasīšanu un rakstīšanu.

BaseExtractor: Pārbauda vai *CorExExtractor* klase satur projekta darbībai nepieciešamās funkcijas.

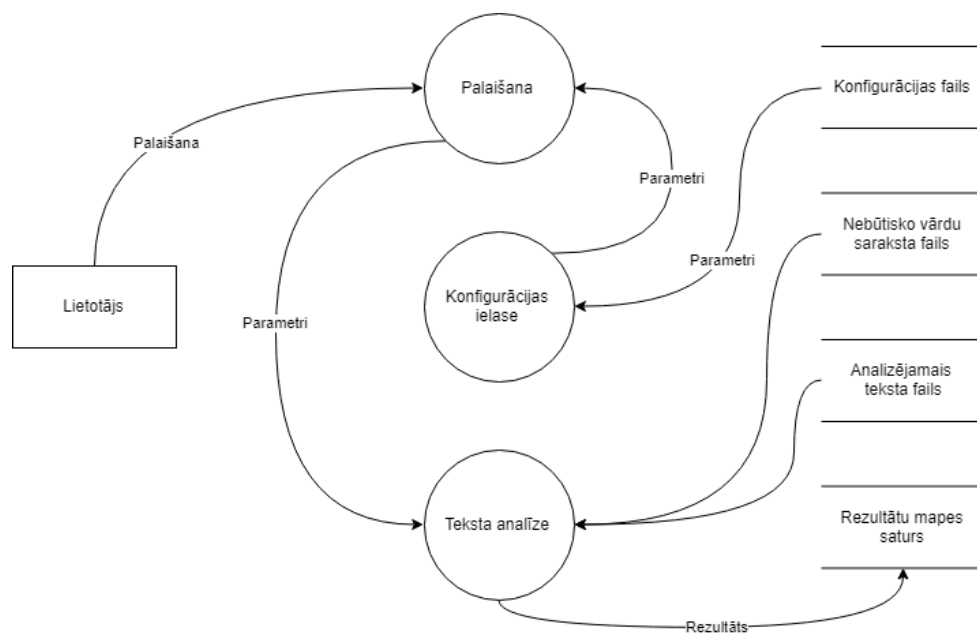
Factory: Izmantojot klases *RunConfiguration* un *CorExConfiguration*, saglabā lietotāja definējamās mainīgās vērtības, kas atrodamas produkta konfigurācijas failā, kā arī klasei *Trainer* dod piekļuvi klases *CorExExtractor* funkcijām.

RunConfiguration: ielasa konfigurācijas faila sadaļu “*RUN*”

CorExConfiguration: ielasa konfigurācijas faila sadaļu “*COREX*”

CorExExtractor: satur vajadzīgās funkcijas, lai nodrošinātu *corextopic*[3] bibliotēkas implementāciju *MMTA*

Attēlā (att.) redzams *MMTA* darbības koncepts. *Corextopic*[3] bibliotēkas implementācijai jānotiek *MMTA* darbības sadaļā “Teksta analīze”. Lai gan pēc vajadzības tā var arī notikt citās sadaļās



Att 1.2.2.1 MMTA darbības koncepts

1.2.3. Lietotāja raksturiezīmes

MMTA lietotājs ir izstrādātās sistēmas lietotājs. Lietotājs ir tāds, kuram ir pasūtītāja piešķirta license izmantot *MMTA* vai arī pasūtītāja darbinieks, kas izmanto *MMTA* darba vajadzībai. *MMTA*. Sistēmas saskarne tiek realizēta izmantojot.csv un .cfg kā arī vai nu .dockerfile, vai nu .exe, failus, kas nosaka, ka *MMTA* lietotājam ir nepieciešamas labas zināšanas darbā ar datoru kā arī vismaz vidējās zināšanas par mašīnmācīšanos un teksta analīzi.

1.2.4. Vispārēji ierobežojumi

Sistēmas darbība ir iespējama izmantojot Windows un Linux operētājsistēmas. Vai arī sistēmas, kas spēj palaist *windows executables* vai *dockerfile*, jo *MMTA* var tikt kompilēts gan kā *windows executable* fails, gan kā *dockerfile* fails.

1.3. Konkrētās prasības

1.3.1. Vispārīgā informācija

CorEx[7] implementācijas funkcionalitātes izklāsts strukturēts skatoties uz katras funkcijas piederību kādam sistēmas konceptuālajam moduļim. Sistēmas *Corex[7]* implementācija moduļi izšķirami sekojoši:

- Palaišana
 - *TopicModeling* - MMTA sākuma punkts
 - *Trainer* - Ekstraktora vadīšana
 - *Factory* - Saskaņojums ar konfigurācijas ielasītājiem un ekstraktoriem
- Konfigurācijas ielase
 - *RunConfiguration* - Palaišanas konfigurācijas ielasītājs
 - *CorExConfiguration* - Ekstraktora konfigurācijas ielasītājs
- Teksta Analīze
 - *FileHelper* - Failu rakstīšana, lasīšana
 - *CorExExtractor* – Ekstraktors

Pēc pasūtītāja lūguma detalizēti tiks apskatītas tikai *CorEx[7]* konfigurācijas un *CorEx[7]* ekstraktijas moduļi, pārējie moduļi tiks aprakstīti, lai būtu saprotama to pamatdarbība.

Prasību specifikācijas trasējamības realizēšanai, katra funkcija tiek identificēta ar unikālu identifikatoru. Identifikatoru veido funkcijas piederība klasei un funkcijas nosaukums rakstīts ar lielajiem burtiem. Klases un funkcijas nosaukums ir atdalāms ar punktu jeb “.” simbolu. Ja klases nosaukumu veido vairāki vārdi, tad vārdi tiek rakstīti kopā (bez atstarpes). Katru vārdu klases nosaukumā raksta ar lielo pirmo burtu, bet pārējos burtus ar mazajiem burtiem. Savukārt vārdu funkcijas nosaukumā atdala ar apakšsvītru jeb “_” simbolu. Visi vārdi funkcijas nosaukumā tiek rakstīti ar mazajiem burtiem. Identifikatora beigās tiek pierakstītas atvērtošās un aizverošās iekavas jeb “()” simbolu kombinācija. Ja identifikators satur tikai klases nosaukumu, tad tas ir uzskatāms par klases identifikatoru.

Identifikatoru piemēri:

- *CorExExtractor.save_models()*-funkcija;

- *CorExConfiguration.save_as_dictionary()* - funkcija;
- *Trainer* - klase;
- *FileHelper* - klase;
- *FileHelper.read_csv()* - funkcija

Visām sistēmas funkcijām ir sekojošas kopīgas iezīmes :

- Visas sistēmā pielietotās simbolu virknes tiek kodētas *UTF8*.
- Ja nav citādāk minēts, tad ievaddati ir obligāti un citādāk izsauc kļūdas ziņojumu par trūkstošo ievaddatu neesamību.
- *None*, *False* un *True* nav klašu vai funkciju identifikatori, tās ir mainīgo vērtības

1.3.2. Funkcionālās prasības

1.3.2.1. Palaišana

1.3.2.1.1. TopicModeling

1.1 tabula

Tēmu modelēšanas programmas sākumpunkts

Identifikators: TopicModeling
Mērķis
Uzsāk <i>MMTA</i> darbību izmantojot <i>Trainer</i> klasi.

1.2 tabula

Tēmu modelēšanas palaišanas funkcija

Identifikators: TopicModeling.main()
Mērķis
Izveido <i>Trainer</i> klases instanci izmantojot <i>Trainer.init()</i> un izsauc <i>Trainer.start_services()</i>
Ievaddati
<ul style="list-style-type: none"> • Konfigurācijas fails
Izvaddati
<ul style="list-style-type: none"> • Konsolē tiek izvadīts paziņojums "Tēmu Modelēšana ir sākta"

1.3.2.1.2. Trainer

1.3 tabula

Tēmu modelēšanas klase Trainer

Identifikators: Trainer
Mērķis
Izveidot saskarni starp <i>TopicModeling</i> klasi un lietotāja izvēlēto teksta analīzes metodi jeb ekstraktoru. Sistēmas izstrādāšanas laikā tiek apskatīts tikai <i>CorExExtractor</i> ekstraktors.

1.4 tabula

Trainer klases inicializācija

Identifikators: Trainer.init()
Mērķis
Inicializē <i>Factory</i> klases instanci, izsauc funkcijas <i>Factory.create_run_configuration()</i> un <i>Factory.create_extractor_configuration()</i>
Ievaddati
<ul style="list-style-type: none">• Konfigurācijas fails

1.5 tabula

Trainer klases darbības sākumpunkts

Identifikators: Trainer.start_services()
Mērķis
Izsauc <i>Trainer.run_training()</i>
Ievaddati
<ul style="list-style-type: none">• Konfigurācijas fails

1.6 tabula

Ekstraktoru izveides, vadības un konfigurācijas failu ielases sākumpunkts

Identifikators: Trainer.run_training()
Mērķis

Izveido <i>Factory</i> klases instanci izmantojot funkciju <i>Factory.init()</i> un izsauc <i>Factory.create_extractor()</i> , <i>CorExExtractor.train()</i> , <i>CorExExtractor.extract()</i> un <i>CorExExtractor.save_models()</i> ,
Ievaddati
<ul style="list-style-type: none"> • Konfigurācijas fails

1.3.2.1.3. Factory

1.7 tabula

Factory

Identifikators: Factory
Mērķis
Veido saskarni ar konfigurācijas faila sadaļu lasītājiem (piemēram <i>RunConfiguration</i>) un ekstraktoriem (piemēram <i>CorExExtractor</i>). Izmantojot klases <i>RunConfiguration</i> un <i>CorExConfiguration</i> , saglabā lietotāja definējamās mainīgās vērtības, kas atrodamas <i>MMTA</i> konfigurācijas failā, klasei <i>Trainer</i> dod piekļuvi klašu <i>CorExExtractor</i> , <i>RunConfiguration</i> , <i>CorExConfiguration</i> funkcijām.

1.7 tabula

Inicializācija

Identifikators: Factory.init()
Mērķis
Izveidot klases <i>Factory</i> instanci
Ievaddati
<ul style="list-style-type: none"> • Nav
Apstrāde
1. Izveido sarakstu, kas satur informāciju par to kādi faili ir vajadzīgi, lai izveidotu konfigurācijas failā izvēlēto ekstraktoru un ielasītu izvēlēta ekstraktora konfigurācijas parametrus no konfigurācijas faila.
Izvaddati
<ul style="list-style-type: none"> • Ekstraktora parametri (<i>CorExConfiguration.init()</i>)
Kļūdu paziņojumi

<ul style="list-style-type: none"> • nav

1.8 tabula

Palaišanas konfigurācijas ielase

Identifikators: Factory.create_run_configuration()
Mērķis
Izveidot klases <i>RunConfiguration</i> instanci izsaucot <i>RunConfiguration.init()</i>
Ievaddati
<ul style="list-style-type: none"> • Konfigurācijas fails
Apstrāde
1. Nav
Izvaddati
<ul style="list-style-type: none"> • Palaišanas parametri (RUN konfigurācijas sadaļa)
Kļūdu paziņojumi
<ul style="list-style-type: none"> • Nav

1.9 tabula

Ekstraktora konfigurācijas ielase

Identifikators: Factory.create_extractor_configuration()
Mērķis
Izveidot ekstraktora klases (<i>CorExExtractor</i>) instanci, ko lietotāja uzstādījis RUN konfigurācijas sadaļā
Ievaddati
<ul style="list-style-type: none"> • Konfigurācijas fails • Ekstraktora tips (“CorEx”)
Apstrāde
2. Apskatot
Izvaddati
<ul style="list-style-type: none"> • Ekstraktora parametri (<i>CorExConfiguration.init()</i>)
Kļūdu paziņojumi
<ul style="list-style-type: none"> •

Ekstraktora klases izveide

Identifikators: Factory.create_extractor()
Mērķis
Izsauc CorExExtractor.init()
Ievaddati
<ul style="list-style-type: none"> • Ekstraktora tips (šajā dokumentā simbolu virkne “CorEx”) • Ekstraktora parametri (<i>CorExConfiguration.init()</i>)
Apstrāde
1.
Izvaddati
•
Kļūdu paziņojumi
•

1.3.2.2. Konfigurācijas ielase

1.3.2.2.1. RunConfiguration

RunConfiguration

Identifikators: RunConfiguration
Mērķis
Ielasīt lietotāja noteiktās vērtības konfigurācijas failā un darbināt <i>MMTA</i> izmantojot šīs vērtības

Inicializācija. Palaišanas konfigurācijas ielase

Identifikators: RunConfiguration.init()
Mērķis
Inicializē <i>RunConfiguration</i> klases instanci, ielasot, pārbaudot un saglabājot konfigurācijas faila sadaļu RUN

Ievaddati
<ul style="list-style-type: none"> • Konfigurācijas fails
Izvaddati
<ul style="list-style-type: none"> • Nav
Apstrāde
1. Pārbauda, vai konfigurācijas parametru vērtības atbilst pareizajiem datu tipiem
Izvaddati
<ul style="list-style-type: none"> •
Kļūdu paziņojumi
<ul style="list-style-type: none"> •

1.3.2.2.2. CorExConfiguration

1.13 tabula

CorExConfiguration

Identifikators: CorExConfiguration
Mērķis

1.14 tabula

Inicializācija. Ekstraktora konfigurācijas ielase

Identifikators: CorExConfiguration.init()
Mērķis
Ievaddati
<ul style="list-style-type: none"> •
Apstrāde
1.
Izvaddati
<ul style="list-style-type: none"> •
Kļūdu paziņojumi

•

1.3.2.3. Teksta analīze

1.3.2.3.1. CorExExtractor

1.15 tabula

CorExExtractotr

Identifikators: CorExExtractor
Mērķis
Veikt <i>corextopic</i> [3] bibliotēkas implementēšanu <i>MMTA</i> , sekojot projekta koda struktūrai pēc pasūtītāja prasībām.

1.16 tabula

Inicializācija

Identifikators: CorExExtractor.init()
Mērķis
Ievaddati
•
Apstrāde
1.
Izvaddati
•
Kļūdu paziņojumi
•

1.17 tabula

Ekstraktora procesu organizēšana

Identifikators: CorExExtractor.train()
Mērķis

<p>Izsaukt vairākas <i>CorExExtractor</i> klases funkcijas, lai, pirmkārt, pārbaudītu <i>CorExExtractor</i> parametru pareizību (<i>CorExExtractor.check_valid_config()</i>); otrkārt, sagatavotu vai nu vārdu lemmatizāciju, vai nu vārdu sakņu atdalīšanu no vārdiem (<i>CorExExtractor.load_stemmer_or_lemmatizer()</i>); treškārt ielasīt nebūtisko vārdu sarakstu (<i>CorExExtractor.load_stopwords()</i>); ceturtkārt, veikt apstrādājamā teksta pirmsapstrādi (<i>CorExExtractor.preprocess_text()</i>); piektkārt, aprēķināt vārdu (terminu jeb angļu valodā <i>term</i>) biežumu (<i>CorExExtractor.calculate_term_frequency()</i>); sestkārt, izveidotu <i>corex</i>[3] bibliotēkā pieejamo mašīnmācīšanās modeli (<i>CorExExtractor.train_with_corex()</i>).</p>
<p>Ievaddati</p> <ul style="list-style-type: none"> • Ievaddatu mapes atrašanās vieta. Simbolu virkne; • Ievaddatu faila nosaukums. Simbolu virkne; • Ievaddatu atdalošais simbols. Simbolu virkne; • Ievaddatu galvenes esamība. Būla tips; • Ievaddatu tekstu saturošās kolonnas numurs. Vesela skaitliska vērtība.
<p>Apstrāde</p> <ol style="list-style-type: none"> 1. Izsauc funkciju <i>CorExExtractor.check_valid_config()</i> 2. Izsauc funkciju <i>CorExExtractor.load_stemmer_or_lemmatizer()</i>; 3. Izsauc funkciju <i>CorExExtractor.load_stopwords()</i>; 4. Izsauc funkciju <i>CorExExtractor.preprocess_text()</i>; 5. Izsauc funkciju <i>CorExExtractor.calculate_term_frequency()</i>; 6. Izsauc funkciju <i>CorExExtractor.train_with_corex()</i>.
<p>Izvaddati</p> <ul style="list-style-type: none"> • Nav
<p>Kļūdu paziņojumi</p> <ul style="list-style-type: none"> • “Ievaddatu mape X neeksistē”, kur X vietā ievieto lietotāja specificētu Ievaddatu mapes atrašanās vietu, ja šāda mape nav atrodama.

1.18 tabula

Rezultātu ekstrakcija

Identifikators: <i>CorExExtractor.extract()</i>
Mērķis

Pēc funkcijas *CorExExtractor.train()* darbības apstrādāt datus izmantojot *CorExModel.transform()[3]*, lai funkcija *CorExExtractor.save_models()* varētu saglabāt rezultātus failos pareizā formātā.

Ievaddati

- Nav

Apstrāde

1. Izsauc funkciju *CorExModel.transform()[3]*
2. Saglabā *CorExModel.transform()[3]* izvadīto matricu *p_y_given_x[3]* (saīsinājums no: *probability of Y for a given X* [7], nozīme:satur matricu, kur kolonnas atbilst tēmām, rindas atbilst teksta dokumentiem un vērtības ir decimāldaļskaitļi lielumā 0 – 1, kas apzīmē izveidotā modeļa pārliecību par katra dotā dokumenta piederību katrai dotai tēmai, kur 0 nozīmē pilnīgu nepiederību un 1 nozīmē pilnīgu piederību konkrētai tēmai) mainīgajā *docTopicDist*;
3. Iterējot cauri *docTopicDist* sastāda mainīgos:
 - *labeledDocs* – matrica ar trim kolonām. Pirmā kolonna satur dokumentu tekstus, otrā - tēmas numurs, ar lielāko pārliecību dotajam dokumentam, trešā – konkrētās tēmas pārliecības lielums;
 - *multilabelData* – matrica ar neregulāru kolonnu skaitu katrā rindā. Katrai rindai ir vismaz viena kolonna. Pirmā kolonna – satur dokumentu tekstus. Otrā kolonna – tēmas numurs, ar lielāko pārliecību dotajam dokumentam, kas pārsniedz lietotāja uzstādīto *ProbabilityThreshold* (pārliecības sliekšņvērtība) konfigurācijas faila sadaļā COREX
 - *docTopicScores* – vārdnīca (*dictionary*), kuras atslēga (*key*) ir katra dokumenta kārtas numurs, un katras atslēgas vērtība (*value*) ir saraksts (*array*) ar apakšsarakstiem (*nested array*), kuri satur divas vērtības: pirmā - tēmas numurs, otrā – modeļa pārliecība par doto tēmu dotajam dokumentam. Katra saraksta apakšsaraksti ir izkārtoti secībā pēc katra apakšsaraksta otrās vērtības (pārliecības) dilstošā secībā.

Izvaddati

- Konsolē izvada paziņojumu “Notiek rezultātu apstrāde”

Kļūdu paziņojumi

•

1.19 tabula

Rezultātu saglabāšana

Identifikators: CorExExtractor.save_models()
Mērķis
Izmantojot <i>FileHelper</i> klases funkcijas saglabāt ar funkciju <i>CorExExtractor.extract()</i> , izveidotos un aizpildītos mainīgos.
Ievaddati
<ul style="list-style-type: none"> Izvaddatu (rezultātu) mapes atrašanās vieta. Simbolu virkne
Apstrāde
<ol style="list-style-type: none"> Saglabā mainīgo <i>labeledDocs</i> failā ar nosaukumu <i>labeledDocuments.csv</i> izmantojot <i>FileHelper.save_to_csv()</i> Saglabā mainīgo <i>multilabelData</i> failā ar nosaukumu <i>multilabelData.csv</i> izmantojot <i>FileHelper.save_to_csv()</i> Saglabā mainīgo <i>docTopicScores</i> failā ar nosaukumu <i>topicScores.json</i> izmantojot <i>FileHelper.save_json()</i> Saglabā mainīgo <i>topics</i> failā ar nosaukumu <i>topicsByKeywords.csv</i> izmantojot <i>FileHelper.save_to_csv()</i>
Izvaddati
•
Kļūdu paziņojumi
•

1.20 tabula

Konfigurācijas vērtību pārbaude

Identifikators: CorExExtractor.check_config_valid()
Mērķis
Pārbaudīt konfigurācijas failā atrodamos COREX parametrus un to datu tipus.
Ievaddati
<ul style="list-style-type: none"> No konfigurācijas faila ielasīto parametru vārdīca.

<ul style="list-style-type: none"> • <i>CorExExtractor.init()</i> izveidota un cieti aizpildīta vārdnīca, kas satur informāciju par to, kāda tipa katram parametram jābūt, un kāda ir parametra noklusējuma vērtība.
Apstrāde
<ol style="list-style-type: none"> 1. Salīdzina pirmās vārdnīcas vērtību datu tipus ar otrajā vārdnīcā definētajām vērtībām salīdzinot vērtības, kurām ir vienādas atslēgas. 2. Ja kāda atslēga ir otrajā vārdnīcā, bet ne pirmajā, tad izvada kļūdas paziņojumu 1 3. Ja kādas pirmās vārdnīcas tips nav atbilstošs, tad izvada kļūdas paziņojumu 2 4. Ja kādas pirmās vārdnīcas vērtība ir <i>None</i>, tad piešķir otrajā vārdnīcā definēto parametra noklusējuma vērtību
Izvaddati
<ul style="list-style-type: none"> • Ja ir izpildījies 4. apstrādes soļa nosacījums, tad konsolē izvada ziņojumu “Parametram X ir piešķirta noklusējuma vērtība Y”, kur X tiek aizvietots ar parametra nosaukumu, pie kura ir izpildījies 4. apstrādes soļa nosacījums un Y ir šī parametra noklusējuma vērtība.
Kļūdu paziņojumi
<ol style="list-style-type: none"> 1. "Parametru X nevar atrast konfigurācijā", kur X tiek aizvietots ar trūkstošo parametra nosaukumu; 2. “Parametram X jābūt ir vai nu tipa Y, vai nu vērtībai None”, kur X ir tiek aizvietots ar neatbilstošā tipa parametra nosaukumu un Y tiek aizvietots ar pareizo parametra tipa nosaukumu.

1.21 tabula

Pirmsapstrādes ielāde

Identifikators: <code>CorExExtractor.load_stemmer_or_lemmatizer()</code>
Mērķis
Pēc konfigurācijas failā zem sadaļas COREX lietotāja izvēlētā teksta pirmsapstrādes metodes ielādēt un izveidot izvēlētās teksta pirmsapstrādes metodes klases instanci.
Ievaddati
<ul style="list-style-type: none"> • Konfigurācijas parametrs <i>applyStemmingOrLemmatization</i>
Apstrāde

<ol style="list-style-type: none"> 1. Ja <i>applyStemmingOrLemmatization</i> vērtība ir simbolu virkne “stemming” izveidot <i>snowballStemmer</i> klasi, kas ir atvērtā pirmkoda bibliotēka. Ja rodas problēmas ielādējot bibliotēku, izvada kļūdu paziņojumu 1. 2. Ja <i>applyStemmingOrLemmatization</i> vērtība ir simbolu virkne “lemmatization” un izveidot <i>wordNetLemmatizer</i> klasi, kas ir atvērtā pirmkoda bibliotēka. Ja rodas problēmas ielādējot bibliotēku, izvada kļūdu paziņojumu 1. 3. Ja <i>applyStemmingOrLemmatization</i> vērtība ir <i>None</i> tad turpināt darbu.
Izvaddati
<ul style="list-style-type: none"> • Ja izpildījies 3. apstrādes soļa nosacījums, tad izvadīt konsolē ziņojumu “Vārdu celmošana un lemmatizēšana nenotiks, jo <i>applyStemmingOrLemmatization</i> vērtība ir <i>None</i>”
Kļūdu paziņojumi
<ol style="list-style-type: none"> 1. “Nav izdevies ielādēt X bibliotēku. Pārbaudiet ieinstalētās bibliotēkas”, kur X ir trūkstošā bibliotēka.

1.22 tabula

Nebūtisko vārdu ielasīšana

Identifikators: CorExExtractor.load_stopwords()
Mērķis
<p>Izņemt nebūtiskos vārdus no teksta dokumentiem teksta dokumentu datu kopā, ja , konfigurācijas failā zem sadaļas COREX lietotājs ir definējis, .csv faila nosaukumu, kas satur nebūtisko vārdu sarakstu.</p>
Ievaddati
<ul style="list-style-type: none"> • Ievaddatu mapes atrašanās vieta. Simbolu virkne; • Konfigurācijas parametrs <i>removeStopwords</i>. Būla tips • Konfigurācijas parametrs <i>stopwordsFileName</i>. Vai nu <i>None</i>, vai nu simbolu virkne.
Apstrāde
<ol style="list-style-type: none"> 1. Ja <i>removeStopwords</i> vērtība ir <i>True</i>, tad veikt apstrādes 2. soli 2. Ja <i>stopwordsFileName</i> vērtība ir <i>None</i> vai tukša simbolu virkne, tad no atvērtā pirmkoda bibliotēkas <i>nlTK.corpus</i> ielādēt <i>stopwords</i> un saglabāt mainīgajā <i>stopWordList</i> kā <i>set</i> tipa sarakstu. Ja rodas problēmas ielādējot bibliotēku, izvada kļūdu paziņojumu 1.

3. Ja neizpildās 2. apstrādes soļa nosacījums, tad ielādēt no ievaddatu mapes failu, kura nosaukums ir vienāds ar mainīgā <i>stopwordsFileName</i> vērtību. Ja fails nav atrodams, tad izvadīt kļūdu paziņojumu 2.
Izvaddati
•
Kļūdu paziņojumi
<ol style="list-style-type: none"> 1. “Nav izdevies ielādēt X bibliotēku. Pārbaudiet ieinstalētās bibliotēkas”, kur X ir trūkstošā bibliotēka. 2. “Nav izdevies atrast failu X mapē Y.”, kur X ir <i>stopwordsFileName</i> vērtība, bet Y ir ievaddatu mapes vērtība.

1.23 tabula

Teksta pirmsapstrāde

Identifikators: CorExExtractor.preprocess_text()
Mērķis
Tokenizēt ielādēto teksta dokumentu saturu
Ievaddati
<ul style="list-style-type: none"> • Teksta dokumentus saturošā faila atrašanās vieta mainīgajā <i>textData</i>. Simbolu virkne; • Ievaddatu atdalošais simbols mainīgajā <i>dataDelimiter</i>. Simbolu virkne; • Ievaddatu galvenes esamība mainīgajā <i>dataFileHasHeader</i>. Būla tips; • Teksta dokumentu datu kopā tekstu saturošās kolonnas numurs mainīgajā <i>dataTextColumn</i>. Vesela skaitliska vērtība.
Apstrāde
<ol style="list-style-type: none"> 1. Ja <i>textData</i> fails nav atrodams izvadīt kļūdu ziņojumu 1.; 2. Izveidot klases <i>FileHelper</i> instanci un mainīgajā <i>documents</i> ielasīt teksta dokumentus izmantojot <i>FileHelper.read_csv()</i>; 3. Iterējot cauri mainīgā <i>documents</i> rindām, katrā iterācijā izmantojot funkciju <i>CorExExtractor.turn_text_into_tokens()</i> tokenizēt <i>documents</i> rindu un izveodotos tokenus apvienot starp tiem ievietojot atstarpī un pievienot mainīgajam <i>documentsProcessed</i>, kas ir saraksts
Izvaddati

<ul style="list-style-type: none"> • <i>documentsProcessed</i>. Saraksts ar apstrādātiem teksta dokumentiem
Kļūdu paziņojumi
<ol style="list-style-type: none"> 1. “Nav idevies atrast failu X”, kur X ir teksta dokumentus saturošā faila atrašanās vietas vērtība

1.24 tabula

Teksta vektorizēšana

Identifikators: CorExExtractor.calculate_term_frequency()
Mērķis
<p>Aprēķināt terminu sastopamību un saglabāt mainīgajā <i>termFrequency</i>, kā arī izveidot un aizpildīt mainīgo <i>featureNames</i>, kam jāsaturs teksta dokumentu datu kopas unikālo vārdu sarakstu alfabētiskā secībā. Izmantojot atvērtā pirmkoda (<i>open-source</i>) bibliotēku <i>countVectorizer</i> pārveidot teksta dokumentu datu kopu par matricu, kuras kolonnas ir unikālo vārdu alfabētiska secība, un rindas ir teksta dokumentu secība, un vērtības attēlo to, cik reizes dotais vārds ir atrodams dotajā dokumentā.</p>
Ievaddati
<ul style="list-style-type: none"> • Apstrādāta teksta dokumentu datu kopa mainīgajā <i>listDocuments</i>
Apstrāde
<ol style="list-style-type: none"> 1. Izmantojot <i>countVectorizer</i> funkciju <i>fit_transform</i> uz mainīgā <i>listDocuments</i> aizpildīt mainīgo <i>termFrequency</i>; 2. Izmantojot bibliotēkas <i>scipy.sparse</i> funkciju <i>csr_matrix</i> apstrādāt <i>termFrequency</i>; 3. Izmantojot <i>countVectorizer</i> funkciju <i>get_feature_names</i> aizpildīt mainīgo <i>featureNames</i>;
Izvaddati
<ul style="list-style-type: none"> • <i>termFrequency</i>. Matrica; • <i>featureNames</i>. Simbolu virkņu saraksts.
Kļūdu paziņojumi
<ul style="list-style-type: none"> • Nav

Corex modeļa trenēšana

Identifikators: CorExExtractor.train_with_corex()
Mērķis
Izveidot <i>corextopic</i> [3] bibliotēkā pieejamo mašīnmācīšanās modeli un saglabāt to kā klasi ar nosaukumu <i>CorExModel</i> izsaucot funkciju <i>CorExModel.corex()</i> [3]
Ievaddati
<ul style="list-style-type: none"> • Ievaddatu faila nosaukums. Simbolu virkne; • ar <i>CorExExtractor.calculate_term_frequency()</i> izveidotā datne saglabāta mainīgajā <i>termFrequency</i>.; • ar <i>CorExExtractor.calculate_term_frequency()</i> izveidotā datne saglabāta mainīgajā <i>featureNames</i>. Neobligāta vērtība. Pēc noklusējuma - <i>None</i>
Apstrāde
<ol style="list-style-type: none"> 1. Izveido <i>CorExModel</i> klasi (mašīnmācīšanās modeli) izmantojot <i>corextopic</i> bibliotēkas iekšējo funkciju <i>corex()</i>[3] (dokumentā līdzvērtīga funkcijai ar identifikatoru <i>CorExModel.corex()</i>[3]) 2. Izmantojot funkciju <i>CorExModel.fit()</i>[3] veic mašīnmācīšanās modeļa pielāgošanu (<i>fitting</i>) teksta dokumentu datu kopai
Izvaddati
<ul style="list-style-type: none"> • Konsolē izvada “Sākusies CorEx inicializācija un trenēšana”
Kļūdu paziņojumi
<ul style="list-style-type: none"> • Nav

1.3.2.3.2. CorExModel

1.26 tabula

CorExModel

Identifikators: CorExModel
Mērķis
Izveidot <i>corextopic</i> [3] bibliotēkā pieejamo mašīnmācīšanās modeļa, klases instanci un izveidot pieeju šīs bibliotēkas funkcijām. Saglabāt instanci ar nosaukumu <i>CorExModel</i>

1.27 tabula

Inicializācija

Identifikators: CorExModel.corex()
Mērķis
Mašīnmācīšanās modeļa izvede izmantojot <i>corextopic</i> [3] bibliotēku.
Ievaddati
<ul style="list-style-type: none">• Konfigurācijas parametrs <i>numOfTopics</i>,• Konfigurācijas parametrs <i>numOfIter</i>,• Konfigurācijas parametrs <i>randomSeed</i>• Funkcijā <i>CorExExtractor.calculate_term_frequency()</i> izveidotais un aizpildītais mainīgais <i>featureNames</i>
Apstrāde
1. Saktīt atsauci par <i>corextopic</i> [3]
Izvaddati
<ul style="list-style-type: none">• Inicializēts mašīnmācīšanās modelis• Saktīt atsauci par <i>corextopic</i>[3]
Kļūdu paziņojumi
<ul style="list-style-type: none">• Saktīt atsauci par <i>corextopic</i>[3]

Modeļa pielāgošana

Identifikators: CorExModel.fit()
Mērķis

Pielāgot (trenēt) mašīnmācīšanās modeli padotajiem datiem. Vairāk par funkciju lasāms: skatot atsauci par <i>corextopic</i> [3]
Ievaddati
<ul style="list-style-type: none"> • Konfigurācijas COREX sadaļas parametrs <i>numOfTopics</i>; • Konfigurācijas COREX sadaļas parametrs <i>numOfIter</i>; • Konfigurācijas DEVELOPER sadaļas parametrs <i>randomSeed</i>; • Funkcijā <i>CorExExtractor.calculate_term_frequency()</i> iegūtais un aizpildītais mainīgais <i>featureNames</i>
Apstrāde
1. Saktīt atsauci par <i>corextopic</i> [3]
Izvaddati
<ul style="list-style-type: none"> • Trenēts mašīnmācīšanās modelis • Saktīt atsauci par <i>corextopic</i>[3]
Kļūdu paziņojumi
<ul style="list-style-type: none"> • Saktīt atsauci par <i>corextopic</i>[3]

1.28 tabula

Rezultātu iegūšana

Identifikators: CorExModel.transform()
Mērķis
Transformēt padotos datus izmantojot trenēto mašīnmācīšanās modeli atpakaļ saņemot <i>p_y_given_x</i> (skatīt funkciju <i>CorExExtractor.extract()</i>) Vairāk par <i>CorExModel.transform()</i> funkciju lasāms skatot atsauci par <i>corextopic</i> [3].
Ievaddati
<ul style="list-style-type: none"> • Mainīgais <i>termFrequency</i>
Apstrāde
1. Saktīt atsauci par <i>corextopic</i> [3]
Izvaddati
<ul style="list-style-type: none"> • Mainīgais <i>p_y_given_x</i> • Saktīt atsauci par <i>corextopic</i>[3]
Kļūdu paziņojumi

- Saktīt atsauci par *corextopic*[3]

1.3.2.3.3. FileHelper

1.29 tabula

FileHelper

Identifikators: FileHelper
Mērķis
Definē un veic standartizētu dažādu formātu failu ielasīšanu un rakstīšanu

1.30 tabula

Datu kopas rakstīšana .csv failā

Identifikators: FileHelper.save_to_csv()
Mērķis
Saglabā datu kopu .csv formāta failā
Ievaddati
<ul style="list-style-type: none"> • Izvadfaila nosaukums. Simbolu virkne • Datu kopa. Saraksts. • Atdalāmais simbols. Simbols
Apstrāde
1. Izmantojot <i>Python</i> iebūvēto funkciju <i>csv</i> datu kopu ieraksta .csv formāta failā
Izvaddati
<ul style="list-style-type: none"> • Fails .csv formātā, kas satur datu kopas datus.
Kļūdu paziņojumi
<ul style="list-style-type: none"> • “Faila nosaukuma paplašinājums nav pareizs, fails tiek saglabāts ar paplašinājumu .csv”

Datu kopas rakstīšana .pkl failā

Identifikators: FileHelper.save_pickle()
Mērķis
Saglabā datu kopu .pkl formāta failā
Ievaddati
<ul style="list-style-type: none"> • Izvadfaila nosaukums. Simbolu virkne

<ul style="list-style-type: none"> • Datu kopa. Saraksts. • Atdalāmais simbols. Simbols
Apstrāde
1. Izmantojot <i>python</i> iebūvēto funkciju <i>pickle</i> datu kopu ieraksta <i>.pkl</i> formāta failā
Izvaddati
<ul style="list-style-type: none"> • Fails <i>.pkl</i> formātā, kas satur datu kopas datus.
Kļūdu paziņojumi
<ul style="list-style-type: none"> • “Faila nosaukuma paplašinājums nav pareizs, fails tiek saglabāts ar paplašinājumu <i>.pkl</i>”

1.31 tabula

Datu kopas rakstīšana *.json* failā

Identifikators: FileHelper.save_json()
Mērķis
Saglabā datu kopu <i>.json</i> formāta failā
Ievaddati
<ul style="list-style-type: none"> • Izvadfaila nosaukums. Simbolu virkne • Datu kopa. Saraksts. • Atdalāmais simbols. Simbols
Apstrāde
1. Izmantojot <i>python</i> iebūvēto funkciju <i>json</i> datu kopu ieraksta <i>.json</i> formāta failā
Izvaddati
<ul style="list-style-type: none"> • Fails <i>.json</i> formātā, kas satur datu kopas datus.
Kļūdu paziņojumi
<ul style="list-style-type: none"> • “Datu formāts neatbilst <i>.json</i> formātam” • Faila nosaukuma paplašinājums nav pareizs, fails tiek saglabāts ar paplašinājumu <i>.json</i>”

1.32 tabula

Datu kopas lasīšana no *.csv* faila

Identifikators: FileHelper.read_csv()
Mērķis
Ielasīt <i>.csv</i> formāta failu datu kopas mainīgajā

Ievaddati
<ul style="list-style-type: none"> • Ievadfaila nosaukums. Simbolu virkne • Atdalāmais simbols. Simbols.
Apstrāde
2. Ielasa .csv formāta failu mainīgajā <i>fileContents</i>
Izvaddati
<ul style="list-style-type: none"> • Mainīgais <i>fileContents</i>. Saraksts.
Kļūdu paziņojumi
<ul style="list-style-type: none"> • “Faila X šifrējums nav lasāms”, kur X ir ievadfaila nosaukums • “Fails X nav atrodams”, kur X ir ievadfaila nosaukums

1.3.2.3.4. Base Extractor

1.33 tabula

BaseExtractor

Identifikators: BaseExtractor
Mērķis
<p>Pārbaudīt, vai klase <i>CorExExtractor</i> satur obligāti vajadzīgās funkcijas. Sekojošām funkcijām ir obligāti jābūt <i>CorExExtractor</i> klasē:</p> <ul style="list-style-type: none"> • <i>init()</i> identifikators: <i>CorExExtractor.init()</i> • <i>load_models()</i> identifikators: <i>CorExExtractor.load_models()</i> • <i>save_models()</i> identifikators: <i>CorExExtractor.save_models()</i> • <i>train()</i> identifikators: <i>CorExExtractor.train()</i> • <i>extract()</i> identifikators: <i>CorExExtractor.extract()</i>

1.4. Ārējā saskarne

1.4.1. Lietotāju saskarne

Lietotāja saskarne ar *MMTA* tiek realizēta izmantojot teksta redaktorus un ar *Windows executable (.exe)* faila palaišanu. Saskarne tiek organizēta sekojoši:

- *MMTA* mape – satur visus vajadzīgos failus, lai varētu sākt veikt teksta analīzi.

- Fails *configuration.cfg* (konfigurācijas fails) - Konfigurācijas failā lietotājs uzstāda izvēlētos parametrus:
 - Ievaddatu mapes atrašanās vieta;
 - Ievaddatu teksta dokumentus saturošā faila nosaukums;
 - Ievaddatu atdalošais simbols;
 - Ievaddatu galvenes esamība;
 - Ievaddatu tekstu saturošās kolonnas numurs;
 - Rezultātu mapes atrašanās vieta;
 - Tēmu skaits;
 - Iterāciju skaits;
 - Atslēgvārdu skaits;
 - Pārlicības sliekšņvērtība;
 - Unikālo vārdu daudzuma sliekšņvērtība;
 - Izvēle starp vārdu celmošanu, lemmatizēšanu vai neapstrādāšanu;
 - Izvēle starp nebūtisku vārdu noņemšanu vai atstāšanu;
 - Lietotāja izvēlēta nebūtisko vārdu saraksta ielāde vai nebūtisko vārdu saraksta ielāde no bibliotēkas bibliotēkas *nlk.corpus* ar nosaukumu *stopwords*;
- Ievaddatu mape – Satur failus, ko lietotājs vēlas izmantot kā *MMTA* ievaddatus. Šie faili iekļauj sevī sekojošus failus:
 - Ievaddatu teksta dokumentus saturošs fails. Obligāts
 - Nebūtisko vārdu sarakstu saturošs fails. Neobligāts
- *Windows executable* fails – Palaiz *MMTA* ar konfigurācijas failā esošajiem parametriem, apstrādājot ievaddatu teksta dokumentus saturošo failu un teksta analīzes rezultātus saglabājot rezultātu mapē.
- Rezultātu mape – satur *MMTA* iegūtos rezultātus, kas sastāv no šādiem failiem: *labeledDocuments.csv*, *multilabelData.csv*, *topicScores.json*, *topicsByKeywords.csv*..

1.4.2. Aparatūras saskarne

Izstrādātā sistēma darbojas uz jebkuras sistēmas, kas ir spējīga palaist *.exe* failus un atbalsta *.csv*, *.json*, *.cfg* failus. Lietotājam būtu ieteicams izmantot sistēmu ar pilnekrāna, darbvirsma atbalstu.

1.4.3. Programmatūras saskarne

MMTA ir izveidots izmantojot *Python 3.5.5*. Lai gan *CorEx[7] implementācija* pati par sevi iekļauj visas tai vajadzīgās bibliotēkas šobrīd esošajā versijā, tomēr produkta attīstības laikā ir jāseko līdzi šo bibliotēku attīstībai un izmaiņām. Pasūtītāja lūguma dēļ tiks pieminētas dokumentā jau minētās bibliotēkas un dažas populārākās bibliotēkas mašīnmācīšanās nozarē, kas tiek izmantotas *CorEx[7] implementācijas ietvaros*. *MMTA* izstrādē jālietoto šādas bibliotēkas:

- *Python 3.5.5* iekšējās bibliotēkas:
 - *logging* <https://docs.python.org/3/library/logging.html>
 - *configparser* <https://docs.python.org/3/library/configparser.html>
 - *csv* <https://docs.python.org/3/library/csv.html>
 - *json* <https://docs.python.org/3/library/json.html>
 - *pickle* <https://docs.python.org/3/library/pickle.html>
 - *os* <https://docs.python.org/3/library/os.html>
 - *sys* <https://docs.python.org/3/library/sys.html>
 - *abc* <https://docs.python.org/3/library/abc.html>
 - *string* <https://docs.python.org/3/library/string.html>
- Ārējās *Python 3.5.5* bibliotēkas:
 - *Numpy 1.16.0* <https://www.numpy.org/>
 - *nltk 3.4* <https://www.nltk.org/https://www.nltk.org/api/nltk.tokenize.html>
 - *corextopic 1.0.2* https://github.com/gregversteeg/corex_topic
 - *sklearn 0.20.2* <https://scikit-learn.org/stable/>
 - *scipy 1.2.0* <https://www.scipy.org/>
 - *pandas 0.23.4* <https://pandas.pydata.org/>

1.4.4. Sakaru saskarne

Izstrādātās sistēmas darbības laikā netiek izmantoti sakari ar citām sistēmām, taču nākotnē sistēmas attīstība varētu iekļaut implementētu komunikāciju ar mājaslapas saskarni izmantojot *API* un *HTTP* vai *HTTPS* protokolus

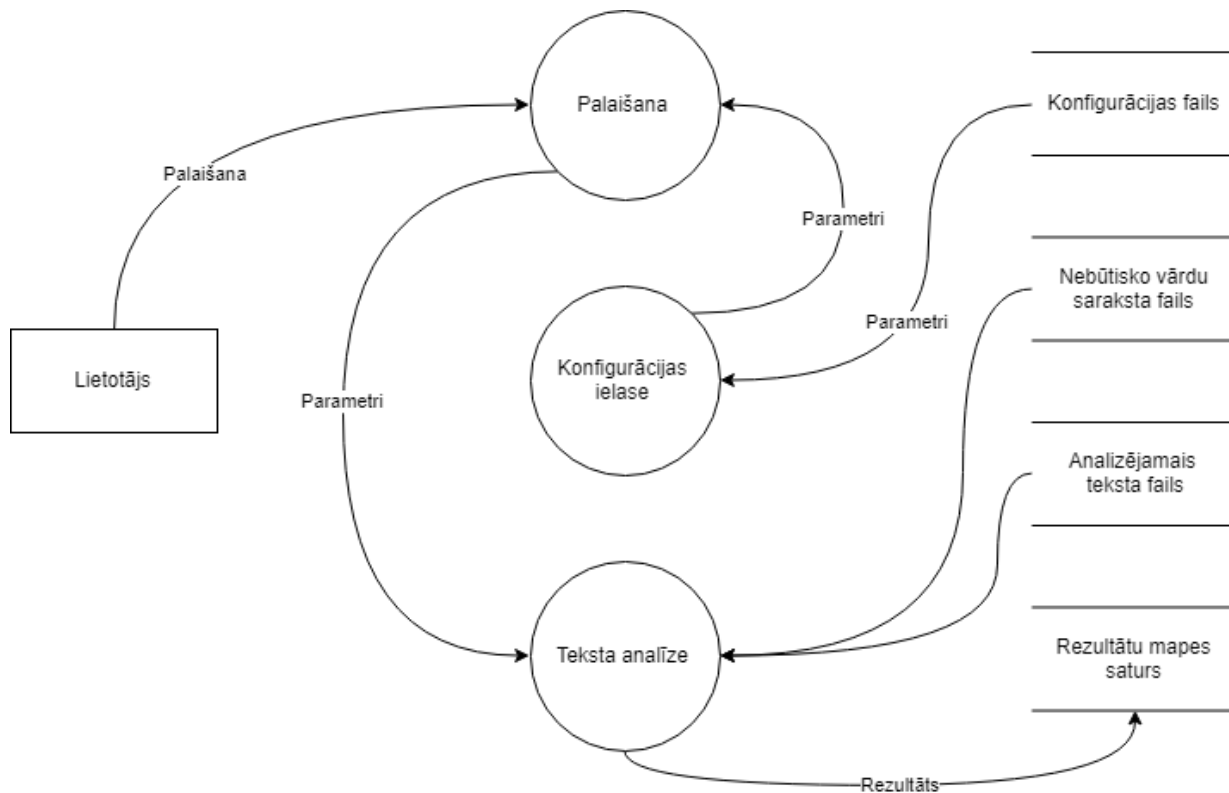
1.5. Nefunkcionālās prasības

1.5.1. Veikstspējas prasības

Apstrādāt 25'000 teksta dokumentu ar 1'000'000 zīmēm 60 sekunžu laikā

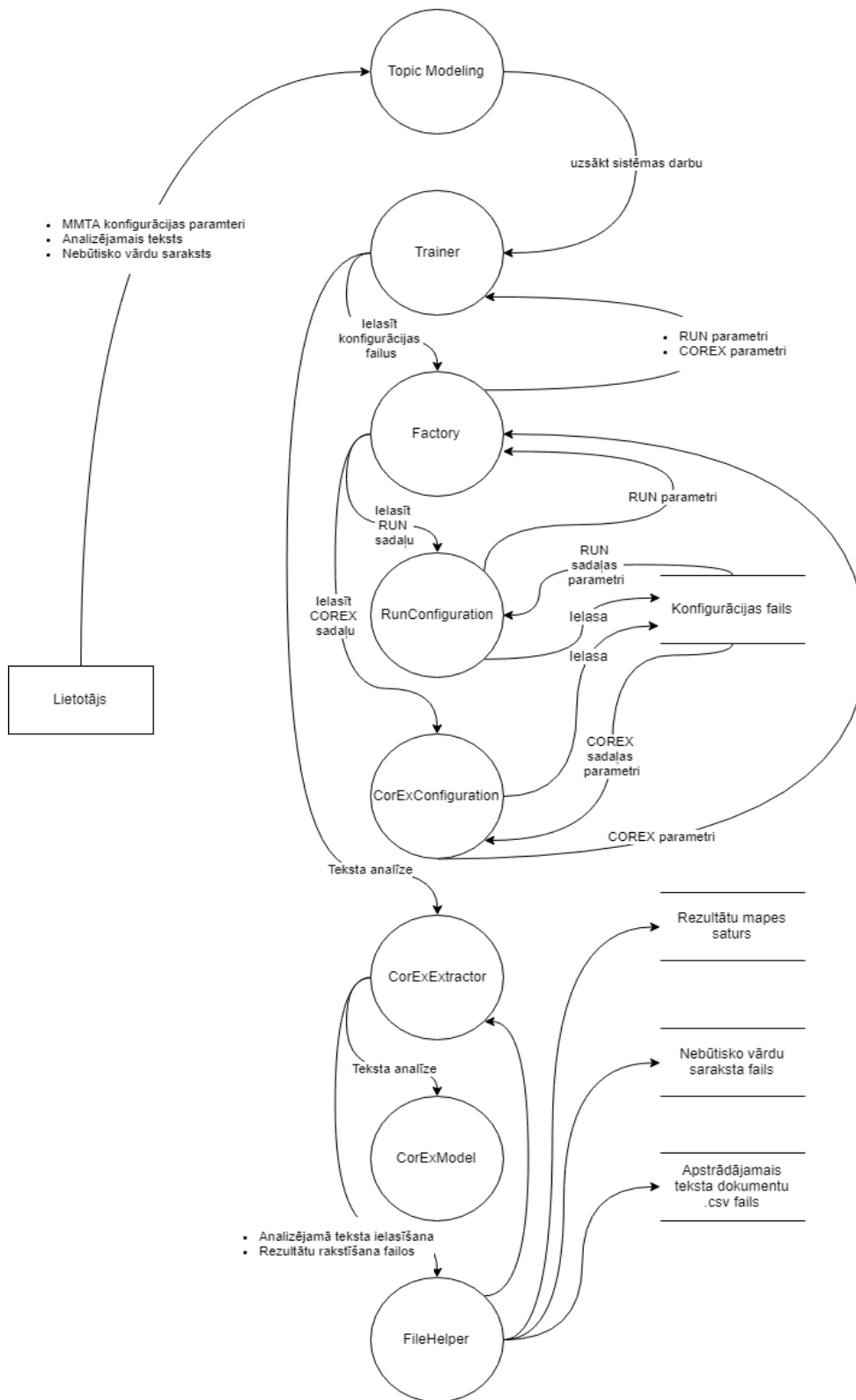
1.6. Konceptuālās datu plūsmas

1.6.1. Sistēmas koncepcija



Att 1.6.1.1 CorEx implementācija sistēmas koncepcija no lietotāja skatpunkta

1.6.2. Pirmā līmeņa Datu plūsmu diagramma –



Att 1.6.2.1 CorEx implementācija sistēmas koncepcija no klašu, moduļu skatpunkta

2. PROGRAMMATŪRAS PROJEKTĒJUMA APRAKSTS

2.1. Ievads

2.1.1. Nolūks

Programmatūras projektējuma apraksts (turpmāk PPA) apraksta kā tehniski tiks izpildīta *corextopic*[3] bibliotēkas implementācija mašīnmācīšanās teksta analizatorā (turpmāk *MMTA*) pēc dokumentā “Programmatūras prasību specifikācija” (turpmāk PPS) uzskaitītajām prasībām. PPA ir paredzēts *MMTA* izstrādātājiem, *corextopic*[3] bibliotēkas implementācijas labojumiem, izprašanai, *MMTA* izstrādātāju iekšējai lietošanai.

2.1.2. Darbības sfēra

Aprakstīts PPS nodaļā “Darbības sfēra”.

2.1.3. Saistība ar citiem dokumentiem

PPA ir lietojams kopā ar PPS

PPA saturs balstīts uz Valsts standarta : LVS 72:1996 "Ieteicamā prakse programmatūras projektējuma aprakstīšanai" prasībām[12].

2.1.4. Pārskats

PPA sastāv no piecām nodaļām:

Ievads – vispārīga informācija par PPA dokumenta saturu

Tehniskie risinājumi – informācija par *corextopic*[3] bibliotēkas implementāciju

Dekompozīcijas apraksts – aprakstītas projektējuma entītijas

Klašu atkarības - aprakstīts kā realizētas implementācijas klašu attiecības

2.2. Tehniskie risinājumi

Sistēma ir izstrādāta izmantojot *Python 3.5.5*, kā arī sekojošās *Python* iekšējās un ārējās bibliotēkas. Ārējā bibliotēkām ir pierakstīts izmantotās bibliotēkas versija. Abu tipu bibliotēkām ir arī minēta tīmekļa vietne, kurā atrodama konkrētās bibliotēkas dokumentācija.

- *Python 3.5.5* iekšējās bibliotēkas:
 - *logging* <https://docs.python.org/3/library/logging.html>
 - *configparser* <https://docs.python.org/3/library/configparser.html>
 - *csv* <https://docs.python.org/3/library/csv.html>
 - *json* <https://docs.python.org/3/library/json.html>
 - *pickle* <https://docs.python.org/3/library/pickle.html>
 - *os* <https://docs.python.org/3/library/os.html>
 - *sys* <https://docs.python.org/3/library/sys.html>
 - *abc* <https://docs.python.org/3/library/abc.html>
 - *string* <https://docs.python.org/3/library/string.html>
- Ārējās *Python 3.5.5* bibliotēkas:
 - *Numpy 1.16.0* <https://www.numpy.org/>
 - *nltk 3.4* <https://www.nltk.org/https://www.nltk.org/api/nltk.tokenize.html>
 - *corextopic 1.0.2* https://github.com/gregversteeg/corex_topic
 - *sklearn 0.20.2* <https://scikit-learn.org/stable/>
 - *scipy 1.2.0* <https://www.scipy.org/>
 - *pandas 0.23.4* <https://pandas.pydata.org/>

2.3. Dekompozīcijas apraksts

2.3.1. Objektu dekompozīcija

PPS prasības tiek dalītas izmantojot objektorientētu dekompozīciju. Katra izstrādes procesā izveidotā klase tiek uzskatīta par objektu. Objekts var saturēt citus objektus, funkcijas, atribūtus, inicializāciju. Objekta inicializācija izveido pašu objektu. Funkcijas ir darbību instrukcijas, kuru rezultātā objekts veic kādu darbību, piemēram, tiek izmainīta kāda atribūta vērtība. Objekts, ko satur cits objekts, ir apakšobjekts. Atribūtam ir nosaukums un tas var saturēt apakšobjektu, skaistlisku vērtībī, simbolu virkni, sarakstu, vārdnīcu un citus datu tipus. Izveidot, atgriezt objekta instanci nozīmē izveidot objektu izsaucot objekta inicializācijas funkciju. Objekts un klase ir savstarpēji aizstājami termini.

Projekta un implementācijas pārskatāmības dēļ projektā sastopamās klases tiks iedalītas trijās kategorijās, pēc klašu nozīmes un darbības mērķa, piederības kādam *MMTA* procesa solim. Katra kategorija satur klases, kas atbilst šī procesa soļa darbībai. Klašu piederība kategorijām ir sekojoša:

- Palaišana
 - *TopicModeling*
 - *Trainer*
- Konfigurācijas ielase
 - *RunConfiguration*
 - *CorExConfiguration* -
- Teksta Analīze
 - *FileHelper* - Failu rakstīšana, lasīšana
 - *CorExExtractor* - Ekstraktors
 - *Factory* - Saskaņotība ar konfigurācijas ielasītājiem un ekstraktoriem

Klases aprakstošās tabulas satur funkciju nosaukumus un to aprakstus, klases atribūtu nosaukumus un atribūtu aprakstus.

Katra klase tiek izveidota atsevišķā failā.

Izstrādātas tika *CorExExtractor* un *CorExConfiguration* klases, kas seko *MMTA* koda struktūrai. Tika veikts darbs arī jau iepriekš izveidotajās klasēs *RunCofiguration*, *Factory*, *Trainer*.

2.3.1.1. Palaišana

Kategorija satur klases, kas atbild par projekta darbības uzsākšanu un darbības organizāciju, izmantojot citu klašu funkcijas.

- Palaišana
 - *TopicModeling* - *MMTA* sākuma punkts
 - *Trainer* - Ekstraktora vadīšana
 - *Factory* - Saskaņojums ar konfigurācijas ielasītājiem un ekstraktoriem
- Konfigurācijas ielase
 - *RunConfiguration* - Palaišanas konfigurācijas ielasītājs
 - *CorExConfiguration* - Ekstraktora konfigurācijas ielasītājs
- Teksta Analīze
 - *FileHelper* - Failu rakstīšana, lasīšana
 - *CorExExtractor* - Ekstraktors

2.3.1.1.1. TopicModeling

2.1 tabula

MMTA galvenais sākumpunkts. Pamatklase.

Atribūti	Apraksts
<i>trainer</i>	Satur <i>Trainer</i> klases instanci.
Funkcijas	Apraksts
<i>main()</i>	Izveido <i>Trainer</i> klases instanci, izsauc <i>Trainer</i> klases funkciju <i>start_service()</i>

2.3.1.1.2. Trainer

2.2 tabula

Ekstraktora vadīšana

Atribūti	Apraksts
<i>factory</i>	Satur klases <i>Factory</i> instanci

<i>extractor</i>	Satur konfigurācijas failā definētā ekstraktora klases instanci (<i>CorExExtractor</i>)
<i>extractorParameters</i>	Satur sarakstu ar ekstraktora konfigurācijas faila parametriem
Funkcijas	Apraksts
<i>init()</i>	Klases inicializācija Izveido <i>Factory</i> klases instanci un saglabā to atribūtā <i>factory</i> Izsauc atribūta <i>factory</i> sekojošās funkcijas: <ul style="list-style-type: none"> • <i>create_run_configuration()</i> • <i>create_extractor_configuration()</i>
<i>start_service()</i>	Izsauc <i>Trainer</i> klases funkciju <i>run_training()</i>
<i>run_training()</i>	Izsauc atribūta <i>factory</i> funkciju <i>create_extractor()</i> Izsauc atribūta <i>extractor</i> sekojošās funkcijas: <ul style="list-style-type: none"> • <i>train()</i> • <i>extract()</i> • <i>save_models()</i>

2.3.1.1.3. Factory

2.3 tabula

Klase, kas atbild par saskarni ar konfigurācijas ielasītājiem, ekstraktoriem

Atribūti	Apraksts
<i>modelTypes</i>	Satur informāciju par to, kā saucas katra pieejamā ekstraktora klase un ekstraktora konfigurācijas ielsīšanas klase.
Funkcijas	Apraksts
<i>init()</i>	Klases inicializācija.
<i>create_run_configuration()</i>	Atgriež <i>RunConfiguration</i> klases instanci
<i>create_extractor_configuration()</i>	Atgriež konfigurācijas failā definētā ekstraktora konfigurācijas ielasītāju (<i>CorExConfiguration</i>)
<i>create_extractor()</i>	Atgriež konfigurācijas failā definētā ekstraktora klases instanci (<i>CorExExtractor</i>)

2.3.1.2. Konfigurācijas ielase

Klases, kas atbild par konfigurācijas faila sadaļu ielasīšanu

2.3.1.2.1. RunConfiguration

2.4 tabula

Klase palaišanas konfigurācijas ielasīšanai.

Atribūti	Apraksts
<i>topicExtractory</i>	Konfigurācijas failā definētais ekstraktora tips (simbolu virkne)
<i>resultsFolder</i>	Konfigurācijas failā definētā izvaddatu (rezultātu) mape (simbolu virkne)
<i>dataFolder</i>	Konfigurācijas failā definētā ievaddatu mape (simbolu virkne)
<i>dataFileName</i>	Konfigurācijas failā definētais ievaddatu (teksta dokumentu kopas) fails (simbolu virkne)
<i>dataDelimiter</i>	Konfigurācijas failā definētā ievaddatu (teksta dokumentu kopas) faila atdalošais simbols (simbols)
<i>dataFileHasHeader</i>	Galvenes esamība konfigurācijas failā definētajā ievaddatu (teksta dokumentu kopas) failam.
<i>dataTextColumn</i>	Tekstu saturošās kolonnas numurs (sākot ar 0) konfigurācijas failā definētajā ievaddatu (teksta dokumentu kopas) failam.
Funkcijas	Apraksts
<i>init()</i>	Klases inicializācija. Ielasa no konfigurācijas faila attiecīgās vērtības un saglabā tās kā klases atribūtus.

2.3.1.2.2. CorExConfiguration

2.5 tabula

Izstrādātā klase ekstraktora konfigurācijas ielasīšanai.

Atribūti	Apraksts
<i>numOfTopics</i>	Konfigurācijas failā definētais tēmu skaits (vesels skaitlis)

<i>numOfIter</i>	Konfigurācijas failā definētais iterāciju skaits (vesels skaitlis)
<i>numOfKeywords</i>	Konfigurācijas failā definētais atslēgvārdu skaits izejas failam <i>topicsByKeywords.csv</i> (vesels skaitlis)
<i>probabilityThreshold</i>	Konfigurācijas failā definētā pārlicības sliekšņvērtība izejas failam <i>multilabelData.csv</i> (skaitlis no 0-1)
<i>maxFeatures</i>	Konfigurācijas failā definētais maksimālais iezīmju skaits pielietotā vektorizētāja (<i>CountVectorizer[]</i>) izveidē (vesels skaitlis)
<i>applyStemmingOrLemmatization</i>	Konfigurācijas failā definētā izvēle starp vārdu apstrādi. Nekāda, vārdu celmošana vai vārdu lemmatizēšana. (<i>None</i> , <i>stemming</i> vai <i>lemmatization</i>)
<i>removeStopwords</i>	Konfigurācijas failā definētā izvēlne nebūtisko vārdu izņemšanas procesam. Neveikt vai veikt. (<i>True</i> vai <i>False</i>)
<i>stopwordsFileName</i>	Konfigurācijas failā definētā ievaddatu mapē atrodamā nebūtisko vārdu sarakstu saturošā faila nosaukums. (<i>None</i> vai simbolu virkne)
Funkcijas	Apraksts
<i>init()</i>	Klases inicializācija. Ielasa no konfigurācijas faila attiecīgās vērtības un saglabā tās kā klases atribūtus.

2.3.1.3. Teksta analīze

Klases, kas atbild par tiešu ievaddatu un izvaddatu apstrādi.

2.3.1.3.1. CorExExtractor

2.6 tabula

Ekstraktora klase

Atribūti	Apraksts
<i>config</i>	Satur vārdnīcu ar ekstraktora parametriem
<i>docTopicDist</i>	Satur <i>CorExModel.transform()</i> izvadīto matricu <i>p_y_given_x[3]</i> (saīsinājums no: <i>probability of Y for a given X[7]</i> , nozīme:satur matricu, kur kolonnas atbilst tēmām, rindas

	atbilst teksta dokumentiem un vērtības ir decimāldaļskaitļi lielumā 0 – 1, kas apzīmē izveidotā modeļa pārliecību par katra dotā dokumenta piederību katrai dotai tēmai, kur 0 nozīmē pilnīgu nepiederību un 1 nozīmē pilnīgu piederību konkrētai tēmai) mainīgajā <i>docTopicDist</i> ;
<i>topics</i>	Satur <i>CorExExtractor</i> klases funkcijas <i>get_topic_results()</i> atgriezto sarakstu
<i>labeledDocuments</i>	Satur <i>CorExExtractor</i> klases funkcijas <i>extract()</i> izveidoto divdimensionālu sarakstu ar trim kolonām: <ol style="list-style-type: none"> 1. teksta dokumenti, 2. tēmas numurs, ar lielāko pārliecību dotajam dokumentam, 3. konkrētās tēmas pārliecības lielums;
<i>multilabelData</i>	Satur <i>CorExExtracotr</i> klases funkcijas <i>extract()</i> izveidoto divdimensionālu sarakstu ar mainīgu kolonu skaitu katrā rindā: <ol style="list-style-type: none"> 1. teksta dokumenti, 2. tēmas numurs, ar lielāko pārliecību dotajam dokumentam, kas pārsniedz lietotāja uzstādīto <i>probabilityThreshold</i> 3. tēmas numurs, ar nākamo lielāko pārliecību dotajam dokumentam, kas pārsniedz lietotāja uzstādīto <i>probabilityThreshold</i>. <p>Kolonnu skaits pēc pirmās kolonnas atkarīgs no tā, cik dotās rindas dokumentam katras tēmas pārliecības vērtība pārsniedz <i>probabilityThreshold</i>. Tik, cik pārsniedz, tik kolonnas pēc pirmās tiek pievienotas dotajai rindai Ja neviena tēmas pārliecības vērtība nepārsniedz <i>probabilityThreshold</i>, tad rinda satur tikai pirmo kolonu.</p>
<i>stemmer</i>	Var saturēt vārdu celmotāja SnowballStemmer instanci no bibliotēkas nltk
<i>lemmatizer</i>	Var saturēt vārdu lemmatizētāja WordNetLemmatizer instanci no bibliotēkas nltk
<i>stopWordList</i>	var saturēt nebūtisko vārdu sarakstu

<i>termFrequency</i>	Satur vektorizētus teksta dokumentus
<i>countVectorizer</i>	Satur vektorizētāja <i>CountVectorizer</i> [9] instanci
<i>corExModel</i>	Satur <i>corextopic</i> [3] instanci
<i>featureNames</i>	Satur teksta dokumentu datu kopas unikālo vārdu sarakstu alfabētiskā secībā.
<i>modelFileNames</i>	Satur sarakstu ar failu nosaukumiem.
<i>CONFIG_TYPES</i>	Satur informāciju par parametru obligāto datu tipu un to noklusējuma vērtībām.
Funkcijas	Apraksts
<i>init()</i>	<p>Klases inicializācija. Aizpilda sekojošos atribūtus:</p> <ul style="list-style-type: none"> • ar tukšu sarakstu: <i>docTopicDist</i>, <i>topicWordDist</i>, <i>topics</i>, <i>labeledDocs</i>, <i>multilabelData</i>, <i>documentsProcessed</i> • ar <i>None</i>: <i>stemmer</i>, <i>lemmatizer</i>, <i>stopWordList</i>, <i>seedList</i>, <i>termFrequency</i>, <i>countVectorizer</i>, <i>corExModel</i>, <i>featureNames</i> • <i>config</i> atribūtu ar sarakstu, kas aizpildīts ar klases <i>Factory</i> funkcijas <i>create_extractor_configuration()</i> ielasītajām vērtībām no konfigurācijas faila. • <i>modelFileNames</i> satru sekojošu sarakstu: “<i>topicsByKeywords.csv</i>, <i>labeledDocuments.csv</i>, <i>topicScores.json</i>, <i>multilabelData.csv</i>” • <i>CONFIG_TYPES</i> satur šādu vārnīcu: <ul style="list-style-type: none"> ○ <i>numOfTopics</i>: ((int, type(None)), 10), ○ <i>numOfIter</i>: ((int, type(None)), 200), ○ <i>numOfKeywords</i>: ((int, type(None)), 10), ○ <i>probabilityThreshold</i>: ((float, type(None)), 0), ○ <i>maxFeatures</i>: ((int, type(None)), 5000), ○ <i>applyStemmingOrLemmatization</i>: ((str, type(None)), 'None'), ○ <i>removeStopwords</i>: ((bool, type(None)), 'True'), ○ <i>stopwordsFileName</i>: ((str, type(None)), 'None'),

	<p>Kur, ieraksts sarakstā sākas ar mainīgā nosaukumu, seko ar datu tipu, tad alternatīvo datu tipu un tad ar noklusējuma vērtību. Piemēram, “<i>numOfTopics: ((int, type(None)), 10),</i>” nozīmē, ka mainīgajam <i>numOfTopics</i> jābūt veselam skaitlim vai arī <i>None</i>, taču, ja ir <i>None</i>, tad vēlākā punktā <i>numOfTopics</i> atribūtam <i>CorExExtractor</i> klases funkcija <i>train()</i>, izmantojot klases <i>CorExExtractor</i> funkciju <i>check_config_valid()</i>, piešķirs noklusējuma vērtību 10.</p>
<i>train()</i>	<p>Secīgi izsauc klases <i>CorExExtractor</i> šādas funkcijas:</p> <ol style="list-style-type: none"> 1. <i>check_config_valid()</i> 2. <i>load_stemmer_or_lemmatizer()</i> 3. <i>load_stopwords()</i> 4. <i>preprocess_text()</i> 5. <i>calculate_term_frequency()</i> 6. <i>train_with_corex()</i>
<i>check_config_valid()</i>	<p>Veic atribūta konfigurācijas parametru pārbaudi, salīdzinot <i>config</i> vārduņicas vērtības ar <i>CONFIG_TYPES</i> atribūtā definētajām parametru raksturīpašībām</p>
<i>load_stemmer_or_lemmatizer()</i>	<p>Ja parametra <i>applyStemmingOrLemmatization</i> vērtība ir:</p> <ol style="list-style-type: none"> 1. simbolu virkne “stemming”, izveidot <i>snowballStemmer</i> instanci no <i>nlk</i> bibliotēkas. 2. simbolu virkne “lemmatization”, izveidot <i>wordNetLemmatizer</i> instanci no <i>nlk</i> bibliotēkas. 3. <i>None</i>, tad turpināt darbu.
<i>load_stopwords()</i>	<p>Ja <i>removeStopwords</i> vērtība ir:</p> <ol style="list-style-type: none"> 1. <i>False</i>, tad turpināt darbu ar tukšu <i>stopWords</i> sarakstu 2. <i>True</i>, tad atribūtā <i>stopWordList</i> ielasīt sekojošās vērtības: ja nebūtisko vārdu saraksta faila nosaukums, kas definēts konfigurācijas failā ir: <ol style="list-style-type: none"> a) netukša simbolu virkne, tad izveidojot <i>FileHelper</i> klases instanci, un izmantojot tās funkciju <i>read_csv()</i>,

	<p>ielasīt nebūtisko vārdu saraksta faila saturu, kas definēts konfigurācijas failā.</p> <p>b) vērtība <i>None</i> vai tukša simbolu virkne, tad ielasīt nebūtisko vārdu sarakstu no <i>nlk</i> bibliotēkas moduļa <i>stopwords</i>.</p>
<i>preprocess_text()</i>	<ol style="list-style-type: none"> 1. Izveidot klases <i>FileHelper</i> instanci, mainīgajā <i>documents</i> ielasīt teksta dokumentus, izmantojot <i>FileHelper</i> klases funkciju <i>read_csv()</i>; 2. Iterējot cauri mainīgā <i>documents</i> rindām, katrā iterācijā izmantojot funkciju <i>CorExExtractor.turn_text_into_tokens()</i> tokenizēt <i>documents</i> rindu un izveidotos tokenus(vārdus) apvienot starp tiem ievietojot atstarpi (izveidot “teikumu”) un pievienot sarakstam <i>documentsProcessed</i>.
<i>calculate_term_frequency()</i>	<ol style="list-style-type: none"> 1. Izmantojot bibliotēkas <i>CountVectorizer</i> funkciju <i>fit_transform</i> uz atribūtu <i>listDocuments</i>, aizpilda atribūtu <i>termFrequency</i>; 2. Izmantojot bibliotēkas <i>scipy</i> funkciju <i>csr_matrix</i> apstrādā atribūtu <i>termFrequency</i>; 3. Izmantojot <i>countVectorizer</i> funkciju <i>get_feature_names</i> aizpilda atribūtu <i>featureNames</i>;
<i>train_with_corex()</i>	<ol style="list-style-type: none"> 1. Izveido <i>CorExModel</i> klases instanci (mašīnmācīsnās modeli) izmantojot <i>corextopic[3]</i> bibliotēkas iekšējo funkciju <i>corex()[3]</i> (tālāk dokumentā funkcija <i>corextopic[3]</i> funkcijas tiks sauktas par klases <i>CorExModel</i> funkcijām) 2. Izmantojot klases <i>CorExModel</i> funkciju <i>fit()[3]</i> veic mašīnmācīsnās modeļa pielāgošanu (<i>fitting</i>) teksta dokumentu datu kopai
<i>extract()</i>	<ol style="list-style-type: none"> 1. Izsauc klases <i>CorExModel</i> funkciju <i>transform()</i>

	<p>2. Atribūtā <i>docTopicDist</i> aglabā klases <i>CorExModel</i> funkcijas <i>transform()[3]</i> izvadīto matricu <i>p_y_given_x[3]</i> (saīsinājums no: <i>probability of Y for a given X[7]</i>, nozīme:satur matricu, kur kolonnas atbilst tēmām, rindas atbilst teksta dokumentiem un vērtības ir decimāldaļskaitļi lielumā 0 – 1, kas apzīmē izveidotā modeļa pārliecību par katra dotā dokumenta piederību katrai dotai tēmai, kur 0 nozīmē pilnīgu nepiederību (0%) un 1 nozīmē pilnīgu piederību (100%) konkrētai tēmai)</p> <p>3. Iterējot cauri <i>docTopicDist</i> aizpilda atribūtus:</p> <p>a) <i>labeledDocs</i> – matrica ar trim kolonām. Pirmā kolonna satur dokumentu tekstus, otrā - tēmas numurs, ar lielāko pārliecību dotajam dokumentam, trešā – konkrētās tēmas pārliecības lielums;</p> <p>b) <i>multilabelData</i> – matrica ar neregulāru kolonnu skaitu katrā rindā. Katrai rindai ir vismaz viena kolonna. Pirmā kolonna – satur dokumentu tekstus. Otrā kolonna – tēmas numurs, ar lielāko pārliecību dotajam dokumentam, kas pārsniedz lietotāja uzstādīto <i>ProbabilityThreshold</i> (pārliecības sliekšņvērtība) konfigurācijas faila sadaļā COREX</p> <p>c) <i>docTopicScores</i> – vārdnīca (<i>dictionary</i>), kuras atslēga (<i>key</i>) ir katra dokumenta kārtas numurs, un katras atslēgas vērtība (<i>value</i>) ir saraksts (<i>array</i>) ar apakšsarakstiem (<i>nested array</i>), kuri satur divas vērtības: pirmā - tēmas numurs, otrā – modeļa pārliecība par doto tēmu dotajam dokumentam. Katra saraksta apakšsaraksti ir izkārtoti secībā pēc katra apakšsaraksta otrās vērtības (pārliecības) dilstošā secībā.</p>
<i>get_topic_results()</i>	Atgriež klases <i>CorExModel</i> funkcijas <i>get_topics()[3]</i> atgriezto sarakstu.

<i>save_models()</i>	<ol style="list-style-type: none"> 1. Saglabā atribūtu <i>labeledDocs</i> failā ar nosaukumu <i>labeledDocuments.csv</i> izmantojot klases <i>FileHelper</i> funkciju <i>save_to_csv()</i> 2. Saglabā atribūtu <i>multilabelData</i> failā ar nosaukumu <i>multilabelData.csv</i> izmantojot klases <i>FileHelper</i> funkciju <i>save_to_csv()</i> 3. Saglabā mai atribūtu nīgo <i>docTopicScores</i> failā ar nosaukumu <i>topicScores.json</i> izmantojot klases <i>FileHelper</i> funkciju <i>save_json()</i> 4. Saglabā atribūtu <i>topics</i> failā ar nosaukumu <i>topicsByKeywords.csv</i> izmantojot klases <i>FileHelper</i> funkciju <i>save_to_csv()</i>
<i>turn_text_into_tokens()</i>	<ol style="list-style-type: none"> 1. Ievadītajam tekstam noņem punktuāciju (pieturzīmes) 2. Katru dokumentu sadala vārdos (tokenizē) 3. Tokena lielos burtus nomaina uz mazajiem burtiem 4. Ja notiek nebūtisko vārdu izņemšana, tad tokenu izdzēš 5. Ja ir izvēlēta lematizēšana vai celmošana, tad tokenam veic attiecīgo apstrādi 6. Atgriež sarakstu ar tokeniem.

2.3.1.3.2. CorExModel

2.7 tabula

Corextopic bibliotēkas klase.

Atribūti	Apraksts
<i>n_hidden</i>	Vērtība vienāda ar konfigurācijas failā norādīto <i>numOfTopics</i>
<i>max_iter</i>	Vērtība vienāda ar konfigurācijas failā norādīto <i>numOfIter</i>
<i>words</i>	Vērtība vienāda ar klases <i>CorExExtractor</i> atribūtu <i>featureName</i> , izmantots, lai varētu atgriezt dokumentus vektorizēšanas darbības rezultātu.
<i>p_y_given_x</i>	Atribūta nosaukums no: <i>probability of Y for a given X</i> [7], nozīme:satur matricu, kur kolonnas atbilst tēmām, rindas atbilst

	teksta dokumentiem un vērtības ir decimāldaļskaitļi lielumā 0 – 1, kas apzīmē izveidotā modeļa pārliecību par katra dotā dokumenta piederību katrai dotai tēmai, kur 0 nozīmē pilnīgu nepiederību un 1 nozīmē pilnīgu piederību konkrētai tēmai
Funkcijas	Apraksts
<i>corex()</i>	Skatīt <i>corextopic</i> [3] bibliotēkas atsaucē atsauci
<i>fit()</i>	Skatīt <i>corextopic</i> [3] bibliotēkas atsaucē atsauci
<i>transform()</i>	Skatīt <i>corextopic</i> [3] bibliotēkas atsaucē atsauci

2.3.1.3.3. FileHelper

2.8 tabula

Klase, kas atbild par failu rakstīšanu, lasīšanu

Funkcijas	Apraksts
<i>init()</i>	Klases inicializācija
<i>read_csv()</i>	Ielasa .csv formāta failu mainīgajā <i>fileContents</i> un atgriež to.
<i>save_to_csv()</i>	Izmantojot <i>python</i> iebūvēto funkciju <i>csv</i> datu kopu ieraksta .csv formāta failā
<i>save_pickle()</i>	Izmantojot <i>python</i> iebūvēto funkciju <i>pickle</i> datu kopu ieraksta .pkl formāta failā
<i>save_json()</i>	Izmantojot <i>python</i> iebūvēto funkciju <i>json</i> datu kopu ieraksta .json formāta failā

2.3.1.3.4. BaseExtractor

2.9 tabula

Klase, kas atbild par ekstraktora klases sturktūras pareizību

Funkcijas	Apraksts
<i>init()</i>	Klases inicializācija
<i>load_models()</i>	Pārbauda tādas pašas nosaukuma funkcijas esamību <i>CorExExtractor</i> klasē
<i>save_models()</i>	Pārbauda tādas pašas nosaukuma funkcijas esamību <i>CorExExtractor</i> klasē

<i>train()</i>	Pārbauda tādas pašas nosaukuma funkcijas esamību <i>CorExExtracor</i> klasē
<i>extract()</i>	Pārbauda tādas pašas nosaukuma funkcijas esamību <i>CorExExtracor</i> klasē

3. TESTĒŠANAS DOKUMENTĀCIJA

3.1. Automātiskie vienībtesti un funkcionālie testi

Sistēmas *MMTA* izmaiņu pieteikuma gadījumā tiek veikti attālināti automātiski testi izmantojot *Azure DevOps* izveidotos *pipelines* (konveijerapstrādes elements). Taču pirms izstrādātājs veic izmaiņu pieteikumu, tam jāveic ir lokālā vienībtestēšana uz ierīces, uz kuras ir veikta izstrāde, izmantojot *Python 3.5.5* iekšējo vienībtestešanas satvaru (*framework*) ar nosaukumu *unittest*, komandrindā ierakstot *python -m unittest* komandu (<https://docs.python.org/3.5/library/unittest.html>), ar komandrindu atvērtajā mapē tiek atrasti un palaisti visi faili, kā skripti (skat. vārdnīcu), kas satur *unittest* moduli.

```
(acn_malta_topics_build_7) C:\Users\mikus.kalnins\source\repos\Topic-Modeling>python -m unittest
-----
Ran 122 tests in 504.590s

OK
(acn_malta_topics_build_7) C:\Users\mikus.kalnins\source\repos\Topic-Modeling>
```

Att 1.6.1.1 MMTA vienībtestēšana izmantojot komandrindu

3.2. Vienībtestēšanas rezultāti

Pēc *corextopic*[3] bibliotēkas implementācijas tika izveidoti lokālie vienībtesti, kas pārbaudītu *MMTA* sistēmai jaunpievienoto funkcionalitāti, pārbaudot katra faila darbību zemāk redzamajās tabulās. Pēc šo vienībtestu sekmīgas iziešanas, tika palaisti arī pārējie *MMTA* testi, lai pārbaudītu sistēmas integritāti ar *MMTA* sistēmu kopumā, kas sākotnēji nebija sekmīgi, taču pēc izmaiņu veikšanas gan implementācijā, gan testu papildināšana, rezultāti beidzot bija apmierinoši. Zemāk tiks aplūkoti implementācijas pārbaudīšanai izstrādātie testi, un veiktās izmaiņas, papildinājumi pārējos *MMTA* testos, jo pēc pasūtītāja lūguma, ļauts dokumentēt tikai *corextopic*[3] implementācijas ietvaros izstrādātos testus.

Tā kā katrs tests tiek palaists kā skripts un tiek veidots no atsevišķām funkcijām un atsevišķām klasēm ar klases funkcijām un klases atribūtiem, tad katra atsevišķa funkcija un klase tiks aprakstīta atsevišķi. Klases un to funkcijas un atribūti tiks aprakstīti līdzīgi kā PPA 2.3.1 nodaļā “Objektu dekpozīcija” jeb izmantojot objektus. Atsevišķās funkcijas savukārt tiks aprakstītas līdzīgi kā PPS 1.3.2 nodaļā “Funkcionālās prasības”. Klases var izsaukt atsevišķās funkcijas, kas

atrodas dotajā testā, kā arī savrīgi ir pieminēt, ka atsevišķās funkcijas ir kodētas katra testa augšējā daļā jeb pirms klasēm. Tātad katra testa apraksts satur divas daļas:

- Funkcijas – atsevišķās funkcijas, apraksta veids balstīts uz PPS 1.3.2 nodaļu “Funkcionālās prasības”;
- Klases – atsevišķās klases un to apraksta veids balstīts uz PPA 2.3.1 nodaļu “Objektu dekopozīcija”.

Rezultāts: pēc izstrādes veikšanas un vienībtestu palaišanas, viss atbilst testu prasībām.

3.2.1. TestCorExExtractor

Tests pārbauda PPA sadaļā 2.3.1.3.1. “CorExExtractor” aprakstītās klases darbību

3.2.1.1. Funkcijas

3.1 tabula

Set_up_mock_corex_training_data()

Nosaukums: set_up_mock_corex_training_data()
Mērķis
Izveidot ievaddatu teksta dokumentu datu kopas .csv failu ievaddatos norādītajā mapē
Ievaddati
<ul style="list-style-type: none"> • <i>dataFolderPath</i> – mape, kurā saglabāt teksta dokumentu datu kopas .csv failu
Apstrāde
<ol style="list-style-type: none"> 1. Izveidot un aizpildīt mainīgo <i>trainingData</i>, kas ir saraksts, kura apkšsaraksti ir teksta dokumenti; 2. Izmantojot <i>Python 3.5.5</i> bibliotēku <i>csv</i> izveidot mapē <i>dataFolderPath</i> failu ar nosaukumu “<i>data.csv</i>”, kas saturēs <i>trainingData</i> mainīgā saturu.
Izvaddati
<ul style="list-style-type: none"> • nav
Kļūdu paziņojumi
<ul style="list-style-type: none"> • nav

Nosaukums: set_up_mock_corex_models()
Mērķis
Izveidot ievadītajā mapē <i>modelFolder</i> tos failus, kuri ir atzīmēti ar vērtību <i>True</i> , piemēram, ja <i>icludeCv</i> ir <i>True</i> , tad izveidot failu “ <i>corExCv.pkl</i> ”, kas aprakstīts 6. apstrādes punktā.
Ievaddat
<ul style="list-style-type: none"> • <i>modelFolder</i>, simbolu virkne; • <i>includeSettings</i>, Būla datu tips, noklusējuma <i>True</i>; • <i>icludeCv</i>, Būla datu tips, noklusējuma <i>True</i>; • <i>includeTf</i>, Būla datu tips, noklusējuma <i>True</i>; • <i>includeModel</i>, Būla datu tips, noklusējuma <i>True</i>; • <i>includeDocTopicDist</i>, Būla datu tips, noklusējuma <i>True</i>; • <i>includeTopicWordDist</i>, Būla datu tips, noklusējuma <i>True</i>; • <i>includeTopics</i>, Būla datu tips, noklusējuma <i>True</i>; • <i>includeLabeledDocs</i>, Būla datu tips, noklusējuma <i>True</i>.
Apstrāde
<ol style="list-style-type: none"> 1. Ja <i>includeSettings</i>, tad izveido un aizpilda mainīgo <i>corExConfigValues</i> ar testa vārnīcu, kas satur cieti kodētas <i>CorEx</i> ekstraktora konfigurācijas testa vērtības. Saglabā mapē ar nosaukumu, kas vienāds ar <i>modelFolder</i>, failā ar nosaukumu “<i>corExSetting.json</i>”. 2. Izveido un aizpilda mainīgo <i>mockData</i>, kas satur teksta dokumentu kopas sarakstu, 3. Izveido un aizpilda mainīgo <i>mockCvValue</i> (saīsinājums no “<i>mock count vecotrizer value</i>”) ar inicializētu <i>CountVecotrizer</i>[9] instanci, 4. Izveido un aizpilda mainīgo <i>mockTfValues</i> (saīsinājums no “<i>mock term frequency values</i>”) vektorizējot mainīgo <i>mockData</i> ar vektorizētāja instances <i>mockCvValue</i> funkciju <i>fit_transform()</i>[9], 5. ja <i>includeTf</i> (saīsinājums no “<i>include term frequency</i>”) vērtība ir <i>True</i>, tad <i>mockTfValues</i> saturu saglabā mapē ar nosaukumu, kas vienāds ar <i>modelFolder</i>, failā ar nosaukumu “<i>corExTf.pkl</i>”

6. ja *includeCv* (saīsinājums no “*include Count Vectorizer*”) vērtība ir *True*, tad *mockCvValues* saturu saglabā mapē ar nosaukumu, kas vienāds ar *modelFolder*, failā ar nosaukumu “*corExCv.pkl*”
7. ja *includeModel* vērtība ir *True*, tad izveido un aizpilda mainīgo *mockModelValues* izmantojot *corextopic[3]* bibliotēkas funkciju *corex()[3]* *corextopic[3]* klases instances izveidei un funkciju *fit()[3]* izveidotās instances (māšīnmācīšanās modeļa pielāgošana) uz mainīgā *mockTfValues*. Beigās mainīgā *mockModelValues* saturu saglabā mapē ar nosaukumu, kas vienāds ar *modelFolder*, failā ar nosaukumu “*corExModel.pkl*”
8. ja *includeDocTopicDist* (saīsinājums no “*include document - topic distribution*”) vērtība ir *True*, tad izveido un aizpilda mainīgo *docTopicDistValues* ar testa vērtībām un mainīgā *docTopicDistValues* saturu saglabā mapē ar nosaukumu, kas vienāds ar *modelFolder*,” failā ar nosaukumu “*corExDocumentTopicMatrix.pkl*”
9. ja *includeTopicWordDist* (saīsinājums no “*include topic – word distribution*”) vērtība ir *True*, tad izveido un aizpilda mainīgo *topicWordDistValues* ar testa vērtībām un mainīgā *topicWordDistValues* saturu saglabā mapē ar nosaukumu, kas vienāds ar *modelFolder*,” failā ar nosaukumu “*corExTopicWordMatrix.pkl*”
10. ja *includeTopics* vērtība ir *True*, tad izveido un aizpilda mainīgo *topicsByKeywordsValues* ar testa vērtībām un mainīgā *topicsByKeywordsValues* saturu saglabā mapē ar nosaukumu, kas vienāds ar *modelFolder*, failā ar nosaukumu “*corExTopicsByKeywords.pkl*”
11. ja *includeLabeledDocs* (saīsinājums no “*labeled documents*”) vērtība ir *True*, tad izveido un aizpilda mainīgo *labeledDocsValues* ar testa vērtībām un mainīgā *labeledDocsValues* saturu saglabā mapē ar nosaukumu, kas vienāds ar *modelFolder*, failā ar nosaukumu “*corExLabeledDocuments.pkl*”

3.2.1.2. Klases

3.3 tabula

TestCorExExtractorSetUp

Funkcijas	Apraksts
<i>setUp()</i>	neko

<i>test_initialize_corex_extractor()</i>	izveido <i>CorExExtractor</i> instanci ar tukšiem parametriem un pārbauda vai: <ol style="list-style-type: none"> 1. izveidotā ekstraktora tips ir ar nosaukumu “CorEx” 2. pārbauda vai ekstraktorā esoša mainīgā <i>docTopicDist</i> saturs ir tukšs 3. pārbauda vai ekstraktorā esoša mainīgā <i>modelFileNames</i> saturā ir 12 elementi 4. pārbauda vai ekstraktorā esoša mainīgā <i>CONFIG_TYPES</i> saturā ir 11 elementi
<i>tearDown()</i>	neko

3.4 tabula

TestCorExExtractorWithoutExistingModel

Funkcijas	Apraksts
<i>setUp()</i>	<ol style="list-style-type: none"> 1. Inicializē klasi 2. Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i>. 3. Izveido un aizpilda mainīgo <i>trainingData</i>, ko saglabā failā ar nosaukumu “data.csv” mapē <i>tempDir</i>. 4. Izveido un aizpilda mainīgo <i>mockStopwordsData</i>, ko saglabā failā ar nosaukumu “mock_stopwords.csv” mapē <i>tempDir</i>. 5. Izveido mainīgo <i>extractorParams</i> un aizpilda to ar testa vērtībām. 6. Mainīgajā <i>extractor</i> izveido <i>CorExExtractor</i> klases instanci izmantojot <i>extractorParams</i> vērtības.
<i>test_training_without_stemming</i>	Uzstāda <i>extractor.config.applyStemmingOrLemmatization</i> uz <i>None</i> , izmanto <i>extractor.train()</i> , pārbauda vai sakrīt <i>self.extractor.countVectorizer.vocabulary_</i> sakrīt ar <i>expectedVocabulary</i> , kas ir cieti iekodēta šajā funkcijā.

<i>test_training_with_stemming</i>	Līdzīgi kā <i>test_training_without_stemming</i> , bet ar <i>extractor.config.applyStemmingOrLemmatization</i> vērtību uzstādītu uz “ <i>stemming</i> ”
<i>test_turn_text_into_tokens</i>	Pārbauda vai <i>extractor.turn_text_into_tokens</i> funkcija izmet kļūdas paziņojumu uz skaitliskas vērtības, bet atgriež gaidīto rezultātu uz funkcijai padotā dotā teksta.
<i>test_remove_punctuation</i>	Pārbauda vai <i>extractor.test_remove_punctuation</i> funkcija izmet kļūdas paziņojumu uz skaitliskas vērtības, bet atgriež gaidīto rezultātu par funkcijai padoto tekstu.
<i>test_save_models</i>	Pārbauda vai <i>extractor.test_remove_punctuation</i> funkcija izmet kļūdas paziņojumu uz skaitliskas vērtības, bet izveido gaidītos failus, ja fukcijai padota <i>tempDir</i> apvienojums ar “ <i>results</i> ”.
<i>test_load_stemmer</i>	Pārbauda vai ar ar konkrētiem uzstādījumiem tiek izmestas kļūdas paziņojumi, bet ar pareizajiem vai izsaucot <i>extractor.load_stemmer_or_lemmatizer</i> <i>extractor.stemmer</i> satur <i>SnowballStemmer</i> instanci un vai <i>extractor.lemmatizer</i> satur <i>None</i>
<i>test_load_lemmatizer</i>	Pārbauda vai ar ar konkrētiem uzstādījumiem tiek izmestas kļūdas paziņojumi, bet ar pareizajiem vai izsaucot <i>extractor.load_stemmer_or_lemmatizer</i> <i>extractor.lemmatizer</i> satur <i>WordNetLemmatizer</i> instanci un vai <i>extractor.stemmer</i> satur <i>None</i>
<i>test_load_stopwords</i>	<ul style="list-style-type: none"> • Pārbauda vai ar uzstādījumu <i>removeStopwords = False</i> atribūts <i>extractor.stopWordList</i> vienāds ar <i>None</i> • Pārbauda vai ar uzstādījumu <i>removeStopwords = False</i> un <i>stopwordsFileName = None</i> atribūts <i>extractor.stopWordList</i> ir <i>set</i> instance • Pārbauda vai ar uzstādījumu <i>removeStopwords = False</i> un <i>stopwordsFileName = mock_stopwords.csv</i> atribūts <i>extractor.stopWordList</i> ir <i>set</i> datu tipa instance

<i>test_train_with_corex</i>	<ul style="list-style-type: none"> • Pārbauda vai izmet kļūdu paziņojumu, ja trūkst kāds no ievaddatiem izsaucot funkciju <i>extractor.train_with_corex()</i>
<i>test_check_config_valid</i>	<ul style="list-style-type: none"> • Pārbauda vai ar izveidotu nepareiz parametru tiek izmests kļūdu paziņojums • Un vai ar nepareiza tipa parametru vērtībām arī tiek izmests kļūdu paziņojums
<i>tearDown()</i>	Izdzēš <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē

3.5 tabula

TestCorExExtractorWithExistingModel

Funkcijas	Apraksts
<i>setUp()</i>	<ol style="list-style-type: none"> 1. Inicializē klasi 2. Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i>. 3. Izveido mainīgo <i>extractorParams</i> un aizpilda to ar testa vērtībām. 4. Mainīgajā <i>extractor</i> izveido <i>CorExExtractor</i> klases instanci izmantojot <i>extractorParams</i> vērtības. 5. Pagaidu mapē izveido mapi “results” un saglabā mainīgajā <i>resultsFolder</i>; 6. Pagaidu mapē izveido mapi “data” un saglabā mainīgajā <i>dataFolder</i>; 7. Pagaidu mapē izveido mapi “models” un saglabā mainīgajā <i>modelFolder</i> 8. Izsauc testa funkciju <i>set_up_mock_corex_training_data()</i> padodod <i>dataFolder</i> 9. Izsauc testa funkciju <i>set_up_mock_corex_models()</i> padodod <i>modelFolder</i>

	<p>10. Izsauc funkciju <i>extractor.load_models()</i> padodod <i>modelFolder</i> un saglabājot instanci mainīgajā <i>corExExtractor</i>.</p> <p>11. Izveido un aizpilda <i>extractor</i> instances sekojošus mainīgos:</p> <ul style="list-style-type: none"> ○ <i>labeledDocs</i> ar cieti kodētām vērtībām ○ <i>docTopicScores</i> ar cieti kodētām vērtībām ○ <i>documents</i> ar cieti kodētām vērtībām ○ <i>docTopicDist</i> ar cieti kodētām vērtībām ○ <i>topicWordDist</i> ar cieti kodētām vērtībām ○ <i>termFrequency</i> ar mainīgā <i>corExExtractor</i> atribūtu <i>termFrequency</i> ○ <i>corExModel</i> ar mainīgā <i>corExExtractor</i> atribūtu <i>corExModel</i>
<i>test_load_models</i>	<ul style="list-style-type: none"> • Pārbauda vai <i>load_models()</i> uz neeksistējošas mapes izmet kļūdu ziņojumu • Pārbauda, vai netiek izmesti kļūdu ziņojumi uz <i>modelFolder</i>
<i>test_extract</i>	<ul style="list-style-type: none"> • Pārbauda vai <i>extractor</i> instance ir <i>CorExExtracotr</i> instance • Pārbauda vai <i>extractor.extract()</i> neizmet kļūdu paziņojumus
<i>test_train_with_corex</i>	<ul style="list-style-type: none"> • Pārbauda vai tiek imzesta kļūda ar pārlietu liela parametra <i>maxFeatures</i> lielumu izsaucot funkciju <i>train_with_corex()</i> • Pārbauda vai netiek izmesti kļūdu paziņojumi, ja ar parametra <i>maxFeatures</i> lielums ir atiblstošs izsaucot funkciju <i>train_with_corex()</i>
<i>test_save_models</i>	<ul style="list-style-type: none"> • Izmanto <i>CountVectorizer</i> funkciju uz <i>extractor</i> atribūtu <i>documents</i>, izmanto <i>save_models</i> uz <i>resultsFolder</i> un pārbauda vai ir izveidoti vajadzīgie faili.
<i>test_multilabelData_is_list</i>	<ul style="list-style-type: none"> • Pārbauda vai <i>multilabel</i> atribūts ir <i>list</i> tipa

<i>test_multilabelData_list_content</i>	<ul style="list-style-type: none"> Pārbauda vai <i>multilabel</i> atribūta saturs tiek veidots pareizi izsaucot <i>extract()</i> pārbauda un salīdzinot rezultātu ar funkcijā ģenerēto testa sarakstu.
<i>test_multilabelData_csv_content</i>	<ul style="list-style-type: none"> Pārbauda vai <i>multilabel</i> tiek pareizi saglabāts .csv formātā, veicot <i>extract()</i>, <i>save_models()</i>, ielasot izveidoto .csv failu un salīdzinot tā saturu ar funkcijā ģenerētu testa sarakstu.
<i>tearDown()</i>	Izdzēš <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē

3.6 tabula

TestCorExExtractorWithExistingModelWithoutSettings

Funkcijas	Apraksts
<i>setUp()</i>	<ol style="list-style-type: none"> Inicializē klasi Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i>. Izveido mainīgo <i>extractorParams</i> un aizpilda to ar testa vērtībām. Mainīgajā <i>extractor</i> izveido <i>CorExExtractor</i> klases instanci izmantojot <i>extractorParams</i> vērtības. Pagaidu mapē izveido mapi “results” un saglabā mainīgajā <i>resultsFolder</i>; Pagaidu mapē izveido mapi “data” un saglabā mainīgajā <i>dataFolder</i>; Pagaidu mapē izveido mapi “models” un saglabā mainīgajā <i>modelFolder</i> Izsauc testa funkciju <i>set_up_mock_corex_training_data()</i> padodod <i>dataFolder</i> Izsauc testa funkciju <i>set_up_mock_corex_models()</i> padodod <i>modelFolder</i> un <i>includeSettings</i> ar vērtību <i>False</i>
<i>test_load_models()</i>	<ul style="list-style-type: none"> Pārbauda, vai izsaucot <i>load_models()</i> tiek izmest kļūdas paziņojums, ka kāds fails nav atrasts

<i>tearDown()</i>	Izdzēs <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē
-------------------	--

3.7 tabula

TestCorExExtractorWithExistingModelWithoutTF

Funkcijas	Apraksts
<i>setUp()</i>	<ol style="list-style-type: none"> 1. Inicializē klasi 2. Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i>. 3. Izveido mainīgo <i>extractorParams</i> un aizpilda to ar testa vērtībām. 4. Mainīgajā <i>extractor</i> izveido <i>CorExExtractor</i> klases instanci izmantojot <i>extractorParams</i> vērtības. 5. Pagaidu mapē izveido mapi “results” un saglabā mainīgajā <i>resultsFolder</i>; 6. Pagaidu mapē izveido mapi “data” un saglabā mainīgajā <i>dataFolder</i>; 7. Pagaidu mapē izveido mapi “models” un saglabā mainīgajā <i>modelFolder</i> 8. Izsauc testa funkciju <i>set_up_mock_corex_training_data()</i> padodod <i>dataFolder</i> <p>Izsauc testa funkciju <i>set_up_mock_corex_models()</i> padodod <i>modelFolder</i> un <i>includeTf</i> ar vērtību <i>False</i></p>
<i>test_load_models()</i>	Pārbauda, vai izsaucot <i>load_models()</i> tiek izmest kļūdas paziņojums, ka kāds fails nav atrasts
<i>tearDown()</i>	Izdzēs <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē

3.8 tabula

TestCorExExtractorWithExistingModelWithoutCorExModel

Funkcijas	Apraksts
<i>setUp()</i>	1. Inicializē klasi

	<ol style="list-style-type: none"> 2. Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i>. 3. Izveido mainīgo <i>extractorParams</i> un aizpilda to ar testa vērtībām. 4. Mainīgajā <i>extractor</i> izveido <i>CorExExtractor</i> klases instanci izmantojot <i>extractorParams</i> vērtības. 5. Pagaidu mapē izveido mapi “results” un saglabā mainīgajā <i>resultsFolder</i>; 6. Pagaidu mapē izveido mapi “data” un saglabā mainīgajā <i>dataFolder</i>; 7. Pagaidu mapē izveido mapi “models” un saglabā mainīgajā <i>modelFolder</i> 8. Izsauc testa funkciju <i>set_up_mock_corex_training_data()</i> padodod <i>dataFolder</i> <p>Izsauc testa funkciju <i>set_up_mock_corex_models()</i> padodod <i>modelFolder</i> un <i>includeModel</i> ar vērtību <i>False</i></p>
<i>test_load_models()</i>	Pārbauda, vai izsaucot <i>load_models()</i> tiek izmest kļūdas paziņojums, ka kāds fails nav atrasts
<i>tearDown()</i>	Izdzēš <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē

3.9 tabula

TestCorExExtractorWithExistingModelWithoutCV

Funkcijas	Apraksts
<i>setUp()</i>	<ol style="list-style-type: none"> 1. Inicializē klasi 2. Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i>. 3. Izveido mainīgo <i>extractorParams</i> un aizpilda to ar testa vērtībām. 4. Mainīgajā <i>extractor</i> izveido <i>CorExExtractor</i> klases instanci izmantojot <i>extractorParams</i> vērtības.

	<p>5. Pagaidu mapē izveido mapi “results” un saglabā mainīgajā <i>resultsFolder</i>;</p> <p>6. Pagaidu mapē izveido mapi “data” un saglabā mainīgajā <i>dataFolder</i>;</p> <p>7. Pagaidu mapē izveido mapi “models” un saglabā mainīgajā <i>modelFolder</i></p> <p>8. Izsauc testa funkciju <i>set_up_mock_corex_training_data()</i> padodod <i>dataFolder</i></p> <p>Izsauc testa funkciju <i>set_up_mock_corex_models()</i> padodod <i>modelFolder</i> un <i>includeCv</i> ar vērtību <i>False</i></p>
<i>test_load_models()</i>	Pārbauda, vai izsaucot <i>load_models()</i> tiek izmest kļūdas paziņojums, ka kāds fails nav atrasts
<i>tearDown()</i>	Izdzēš <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē

3.10 tabula

TestCorExExtractorWithExistingModelWithoutDocTopicDist

<i>setUp()</i>	<ol style="list-style-type: none"> 1. Inicializē klasi 2. Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i>. 3. Izveido mainīgo <i>extractorParams</i> un aizpilda to ar testa vērtībām. 4. Mainīgajā <i>extractor</i> izveido <i>CorExExtractor</i> klases instanci izmantojot <i>extractorParams</i> vērtības. 5. Pagaidu mapē izveido mapi “results” un saglabā mainīgajā <i>resultsFolder</i>; 6. Pagaidu mapē izveido mapi “data” un saglabā mainīgajā <i>dataFolder</i>; 7. Pagaidu mapē izveido mapi “models” un saglabā mainīgajā <i>modelFolder</i> 8. Izsauc testa funkciju <i>set_up_mock_corex_training_data()</i> padodod <i>dataFolder</i>
----------------	--

	Izsauc testa funkciju <i>set_up_mock_corex_models()</i> padodod <i>modelFolder</i> un <i>includeDocTopicDist</i> ar vērtību <i>False</i>
<i>test_load_models()</i>	Pārbauda, vai izsaucot <i>load_models()</i> tiek izmest kļūdas paziņojums, ka kāds fails nav atrasts
<i>tearDown()</i>	Izdzēš <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē

3.11 tabula

TestCorExExtractorWithExistingModelWithoutTopicWordcDist

Funkcijas	Apraksts
<i>setUp()</i>	<ol style="list-style-type: none"> 1. Inicializē klasi 2. Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i>. 3. Izveido mainīgo <i>extractorParams</i> un aizpilda to ar testa vērtībām. 4. Mainīgajā <i>extractor</i> izveido <i>CorExExtractor</i> klases instanci izmantojot <i>extractorParams</i> vērtības. 5. Pagaidu mapē izveido mapi “results” un saglabā mainīgajā <i>resultsFolder</i>; 6. Pagaidu mapē izveido mapi “data” un saglabā mainīgajā <i>dataFolder</i>; 7. Pagaidu mapē izveido mapi “models” un saglabā mainīgajā <i>modelFolder</i> 8. Izsauc testa funkciju <i>set_up_mock_corex_training_data()</i> padodod <i>dataFolder</i> <p>Izsauc testa funkciju <i>set_up_mock_corex_models()</i> padodod <i>modelFolder</i> un <i>includeTopicWordDist</i> ar vērtību <i>False</i></p>
<i>test_load_models()</i>	Pārbauda, vai izsaucot <i>load_models()</i> tiek izmest kļūdas paziņojums, ka kāds fails nav atrasts
<i>tearDown()</i>	Izdzēš <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē

TestCorExExtractorWithExistingModelWithoutTopics

Funkcijas	Apraksts
<i>setUp()</i>	<ol style="list-style-type: none"> 1. Inicializē klasi 2. Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i>. 3. Izveido mainīgo <i>extractorParams</i> un aizpilda to ar testa vērtībām. 4. Mainīgajā <i>extractor</i> izveido <i>CorExExtractor</i> klases instanci izmantojot <i>extractorParams</i> vērtības. 5. Pagaidu mapē izveido mapi “results” un saglabā mainīgajā <i>resultsFolder</i>; 6. Pagaidu mapē izveido mapi “data” un saglabā mainīgajā <i>dataFolder</i>; 7. Pagaidu mapē izveido mapi “models” un saglabā mainīgajā <i>modelFolder</i> 8. Izsauc testa funkciju <i>set_up_mock_corex_training_data()</i> padodod <i>dataFolder</i> <p>Izsauc testa funkciju <i>set_up_mock_corex_models()</i> padodod <i>modelFolder</i> un <i>includeTopics</i> ar vērtību <i>False</i></p>
<i>test_load_models()</i>	Pārbauda, vai izsaucot <i>load_models()</i> tiek izmest kļūdas paziņojums, ka kāds fails nav atrasts
<i>tearDown()</i>	Izdzēš <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē

3.13 tabula

TestCorExExtractorWithExistingModelWithoutLabeledDocs

Funkcijas	Apraksts
<i>setUp()</i>	<ol style="list-style-type: none"> 1. Inicializē klasi 2. Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i>.

	<p>3. Izveido mainīgo <i>extractorParams</i> un aizpilda to ar testa vērtībām.</p> <p>4. Mainīgajā <i>extractor</i> izveido <i>CorExExtractor</i> klases instanci izmantojot <i>extractorParams</i> vērtības.</p> <p>5. Pagaidu mapē izveido mapi “results” un saglabā mainīgajā <i>resultsFolder</i>;</p> <p>6. Pagaidu mapē izveido mapi “data” un saglabā mainīgajā <i>dataFolder</i>;</p> <p>7. Pagaidu mapē izveido mapi “models” un saglabā mainīgajā <i>modelFolder</i></p> <p>8. Izsauc testa funkciju <i>set_up_mock_corex_training_data()</i> padodod <i>dataFolder</i></p> <p>Izsauc testa funkciju <i>set_up_mock_corex_models()</i> padodod <i>modelFolder</i> un <i>includeLabeledDocs</i> ar vērtību <i>False</i></p>
<i>test_load_models()</i>	Pārbauda, vai izsaucot <i>load_models()</i> tiek izmest kļūdas paziņojums, ka kāds fails nav atrasts
<i>tearDown()</i>	Izdzēš <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē

3.2.2. TestCorExConfiguration

Tests pārbauda PPA sadaļā 2.3.1.2.2. “CorExConfiguration” aprakstītās klases darbību

3.2.2.1. Funkcijas

Atsevišķu funkciju testam nav.

3.2.2.2. Klases

3.14 tabula

TestCorExParameters

Atribūti	Apraksts

Funkcijas	Apraksts
<i>setUp()</i>	<ol style="list-style-type: none"> 1. Inicializē klasi 2. Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i>.
<i>test_number_of_corex_variables()</i>	<ul style="list-style-type: none"> • Izveido un aizpilda mainīgo <i>expectedCorExVariables</i>, kas satur cieti kodētu sarakstu ar parametru nosaukumiem, kurus sagaidām, ka tik ielasīti izmantojot <i>CorExConfiguration</i> klasi. Veic pārbaud vai ielasītais saturs sakrīt ar mainīgā <i>expectedCorExVariables</i> saturu
<i>test_create_corex_configuration_with_incorrect_values()</i>	<ul style="list-style-type: none"> • Izveido konfigurācijas failu, kas satur tikai COREX sadaļu ar nepareiziem datu tipiem un pārbauda, vai ielasot izmantojot <i>CorExConfiguration</i> klasi, netiek paturētas nepareizās vērtības
<i>test_create_corex_configuration_with_correct_values()</i>	Izveido konfigurācijas failu, kas satur tikai COREX sadaļu ar pareiziem datu tipiem un pārbauda, vai ielasot izmantojot <i>CorExConfiguration</i> klasi, tiek paturētas pareizās vērtības
<i>tearDown()</i>	Izdzēš <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē

3.2.3. TestFactory

Tests pārbauda PPA sadaļā 2.3.1.1.3. “Factory” aprakstītās klases darbību

3.2.3.1. Funkcijas

Atsevišķu funkciju testam nav.

3.2.3.2. Klases

3.15 tabula

SetUpFactoryTest

Atribūti	Apraksts

Funkcijas	Apraksts
<i>setUp()</i>	nekas nenotiek
<i>test_initialize_factory()</i>	<ul style="list-style-type: none"> • Inicializē klasi <i>Factory</i> un pārbauda, vai tās atribūts <i>MODEL_TYPES</i> ir ar pareizu saturu salīdzinot pret cieti kodētu vārdnīcu
<i>tearDown()</i>	nekas nenotiek

3.16 tabula

CorExFunctionalFactoryTest

Atribūti	Apraksts
Funkcijas	Apraksts
<i>setUp()</i>	<ol style="list-style-type: none"> 1. Inicializē klasi 2. Izveido pagaidu mapi un saglabā tās atršanās vietu mainīgajā <i>tempDir</i> 3. Izveido un aizpilda mainīgo <i>CONFIG_TYPES</i> 4. Izveido un aizpilda mainīgo <i>corExConfigValues</i> 5. Saglabā <i>tempDir</i> mapē <i>corExConfigValues</i> mainīgā vērtību failā “<i>corExSettings.json</i>” 6. Izveido un aizpilda mainīgo <i>mockTfValues</i> 7. Saglabā <i>tempDir</i> mapē <i>mockTfValues</i> mainīgā vērtību failā “<i>corExTf.pkl</i>” 8. Izveido un aizpilda mainīgo <i>mockCvValue</i> 9. Saglabā <i>tempDir</i> mapē <i>mockCvValue</i> mainīgā vērtību failā “<i>corExCv.pkl</i>” 10. Izveido un aizpilda mainīgo <i>mockModelValues</i> 11. Saglabā <i>tempDir</i> mapē <i>mockModelValues</i> mainīgā vērtību failā “<i>corExModel.pkl</i>” 12. Izveido un aizpilda mainīgo <i>docTopicDistValues</i> 13. Saglabā <i>tempDir</i> mapē <i>docTopicDistValues</i> mainīgā vērtību failā “<i>corExDocumentTopicMatrix.pkl</i>”

	<p>14. Izveido un aizpilda mainīgo <i>topicWordDistValues</i></p> <p>15. Saglabā <i>tempDir</i> mapē <i>topicWordDistValues</i> mainīgā vērtību failā “corExTopicWordMatrix.pkl”</p> <p>16. Izveido un aizpilda mainīgo <i>labeledDocsValues</i></p> <p>17. Saglabā <i>tempDir</i> mapē <i>labeledDocsValues</i> mainīgā vērtību failā “corExLabeledDocuments.pkl”</p> <p>18. Izveido un aizpilda mainīgo <i>topicsByKeywordsValues</i></p> <p>19. Saglabā <i>tempDir</i> mapē <i>topicsByKeywordsValues</i> mainīgā vērtību failā ar nosaukumu “corExTopicsByKeywords.pkl”</p> <p>20. Izveido <i>Factory</i> klases instanci un saglabā to mainīgajā <i>factory</i></p> <p>21. Izveido un aizpilda mainīgo <i>configFile</i> ar “./TopicModeling.cfg”</p>
<i>test_create_run_configuration()</i>	<ul style="list-style-type: none"> • Izveido un aizpilda mainīgo <i>runParams</i> izmantojot <i>factory</i> instances funkciju <i>create_run_configuration</i> uz faila <i>configFile</i> un pārbauda vai <i>runParams</i> ir <i>RunConfiguration</i> instance
<i>test_create_extractor_configuration()</i>	<ul style="list-style-type: none"> • Izmantojot <i>factory</i> instances funkciju <i>create_extractor_configuration()</i> ar nepareizām vērtībām, pārbauda, vai tiks izsaukts kļūdu paziņojums • Izveido un aizpilda mainīgo <i>corExParams</i> izmantojot <i>factory</i> instances funkciju <i>create_extractor_configuration()</i> padodod failu <i>configFile</i> un ar simbolu virkni “CorEx” pārbauda vai <i>corExParams</i> ir <i>CorExConfiguration</i> instance
<i>test_create_extractor()</i>	<ul style="list-style-type: none"> • Izmantojot <i>factory</i> instances funkciju <i>create_extractor()</i> ar nepareizām vērtībām, pārbauda, vai tiks izsaukts kļūdu paziņojums

	<ul style="list-style-type: none"> • Izveido un aizpilda mainīgo <i>corExParams</i> izmantojot <i>factory</i> instances funkciju <i>create_extractor()</i> padodod failu <i>configFile</i> un ar simbolu virkni "CorEx" un izveido un aizpilda mainīgo <i>corExExtractor</i> izmantojot <i>factory</i> instances funkciju <i>create_extractor()</i> padodod tai simbolu virkni "CorEx" un <i>corExParams</i>. Beigās pārbauda vai <i>corExExtractor</i> ir <i>CorExExtractor</i> instance
<i>tearDown()</i>	Izdzēš <i>tempDir</i> mapi un visus failus un datus, kas ir bijuši saglabāti <i>tempDir</i> mapē

4. PROJEKTA ORGANIZĀCIJA

Papildu funkciju izstrāde mašīnmācīšanās teksta analizatoram (turpmāk *MMTA*) tika organizēta pēc spējās izstrādes apakškategorijas *Scrum* metodoloģijas. Uzsākot projektu tika veiktas sarunas ar pasūtītāju, lai noskaidrotu kādai papildus funkcionalitātei jātiek izstrādātai. Kad tika noskaidrots, ka jāveic ir *corextopic*[3] bibliotēkas implementācija, tad tika veikta iepazīšanās ar *MMTA* sistēmu kopumā un apskatīta un izprasta iepriekšējo *MMTA* izstrādātāju pieeja citu bibliotēku implementācijai. *MMTA* izprašanas laikā tika veikti mazi uzdevumi, un, ar katru izpildīto uzdevu, tika izpildītas pasūtīta prasības, kā arī uzstādītas jaunas prasības un jauni uzdevumi. Pēc katra uzdevuma izpildes tuvojoties pilnīgākai projekta izpildei.

Izstrādes process tika sadalīts sākumā divu, pēc tam triju nedēļu garos intervālos, ko sauc par sprintiem, katra sprinta sākuma tika veikta sprinta plānošana un katra sprinta beigās tika veikts sprinta atskats (*retrospective*) kopā ar iesaistītajām pusēm un komandas biedriem. Katras dienas sākumā tika veiktas *Scrum* tikšanās ar komandas biedriem, kuras laikā katrs komandas biedrs izstāstīja padarīto, sastaptās grūtības vai kādus citus jaunumus, kas saistīti ar *MMTA*. Visas tikšanās notika ar video zvana palīdzību izmantojot *Microsoft Teams* lietotni.

Izstrādes pārvaldība notika izmantojot *Azure DevOps* lietotni, kas uzlaboja darba organizāciju un pārskati. Uzsākot izstrādi, tika noteikts, ka programmaprodukta koda izstrādes termiņš ir 2019. gada 17. maijs, pēc kā jāseko dokumentēšanas procesam līdz 2019.gada 25. maijam.

Patērētā laika un līdzekļu mazināšanai tika izmantoti atvērtā pirmkoda rīki, par kuriem lasāms PPA nodaļā “*Tehniskie risinājumi*”.

Izstrāde veikta izstrādes vidēs *Microsoft Visual Studio Professional 2017 15.9.12*, *Jupyter Notebook 4.4.0* un *Sublime Text Version 3.2, Build 3200*. Versiju kontrole un izstrādāta projekta augšupielāde notika izmantojot *Azure DevOps* un *git-gui version 0.21.GITGUI* ar *git version 2.20.1.windows.1*.

Projekta dokumentēšana tika veikta izmantojot *Microsoft Office 365* lietotni *Microsoft Word* un interneta lietotni *Drawio*[8].

Projekta dokumentēšanas, plānošanas un organizēšanas laikā ņemts vērā rakstūtais LU Datorikas fakultātes e-studiju kursa materiālos no kursiem “*Kvalifikācijas darbs*” un

“Programminženierija”. Darba noformēšanā ievērots rīkojums Nr. 1/69 “Prasības noslīguma darbu (bakalaura, maģistra darbu, diplomdarbu un kvalifikācijas darbu) izstrādāšanai un aizstāvēšanai Latvijas Universitātē”.

The screenshot shows a Jira Product Backlog item for '20736 NLP: Implement CorEx clustering into Topic-Modeling' by Mikus Kalnins. It features a 'Related Work' section with a table of linked items.

Link	State	Latest Update	Comment
22605 code extract()	Done	Updated 5/6/2019	
22601 code init()	Done	Updated 4/25/2019	
22604 code train()	Done	Updated 5/3/2019	
22721 Create a minimal ~viable "product"	Done	Updated 5/6/2019	
22613 Document	In Progress	Updated 5/16/2019	
23007 Fix probability in labeledDocuments	Done	Updated 5/3/2019	
23884 Implement Anchoring	To Do	Updated 8 minutes ago	
22650 Implement the option to use CorExExtractor.py	Done	Updated 4/11/2019	
22608 Investigate attached Jupyter Notebooks	Done	Updated 4/10/2019	
22600 Investigate LDAExtractor.py	Done	Updated 4/10/2019	
23417 MultilabelData.csv	Done	Updated 5/7/2019	
23363 Pull dev branch locally	Done	Updated 5/6/2019	
23533 Pull Request, Merge	In Progress	Updated 8 minutes ago	
22720 Reworking the jupyter Notebook into a POC	Done	Updated 4/11/2019	
23504 Test Executables	Done	Updated 8 minutes ago	

Att 4.1 Uzdevumu pārvalde izmantojot Azure DevOps platformu.

5. KONFIGURĀCIJU PĀRVALDĪBA

CorEx[7] implementācijas konfigurāciju pārvaldība tika īstenota ar versiju vadību. Tādā veidā dodot iespēju sekot līdzi izmaiņām programmaprodukta kodā un dokumentācijā: Šāda veida konfigurāciju pārvaldība ļāva brīvi mainīt sistēmas versijas un tās zarus, ja tika atklātas novirzes kādā no implementācijas daļām.

Versiju pārvaldība īstenota ar *Git* repozitoriju, ko pārvalda izmantojot *git-gui* lietotni vai arī izmantojot izstrādes vidē *Microsoft Visual Studio Professional 2017* pieejamo *git* paplašinājumu. Repozitorijs tika saglabāts *Azure DevOps*[2] platformā. Repozitoriju glabāšanas platformas izvēli veicis ir pasūtītājs

Jebkuras izstrādātāja veiktās dokumentācijas vai *MMTA* repozitorija izmaiņas tiek saglabātas pasūtītāja nodrošinātajā izstrādes ierīcē, izmaiņas tiek aprakstītas un kad izmaiņu veikšana ir beigta, tad tās tiek iesūtītas *Git* repozitorijā, un saglabāts *Azure DevOps*[2] platformā, lai izmaiņas varētu tikt notestētas, pārskatītas, pievienotas izstrādes zaram un pieejams citiem izstrādātājiem.

Aprakstītā konfigurāciju pārvaldība ir sevi attaisnojusi un ir visnotaļ efektīva gan viena, gan vairāku izstrādāju gadījumā. Lai gan protams, ja projekts būtu viena cilvēka izstrādāts, tad varētu iztikt bez *Azure DevOps*[2], bet gan izmantot, piemēram, *DropBox*.

6. KVALITĀTES NODROŠINĀŠANA

Kvalitātes nodrošināšana *CorEx*[7] implementācijas laikā tiek nodrošināta izpildot pasūtītāja izvirzītās prasībās un ieviešot prasībās minētās funkcionālītes. Darba kvalitātes novērtēšana notiek *CorEx*[7] implementācijas gala produkta atbilstība pasūtītāja prasībām.

Zemāk redzams saraksts ar būtiskākajām kvalitātes nodrošināšanas aktivitātēm, kuras tika izpildītas darba veikšanas laikā. Daļa no aktivitātēm ir aprakstītas sīkāk šajā dokumentā. Saraksts nav uzskatāms par galīgu.

- Pasūtītāja *MMTA* sistēmas apguve
- Prasību definēšana (skatīt “*Programmatūras prasību specifikācija*”)
- Plānošana (skatīt “*Projekta organizācija*”)
- Dokumentēšanas izstrāde
- Tehniskā realizācija (skatīt “*Tehniskie risinājumi*”)
- Vienībtestēšana (skatīt “*Testēšanas dokumentācija*”)
- Konfigurāciju pārvaldība (skatīt “*Konfigurāciju pārvaldību*”)
- Konsultācijas ar *MMTA* projekta kolēģiem
- Konsultācijas ar pasūtītāju
- Pasūtītāja nodrošinātas papildus apmācības
- Konsultācijas ar kvalifikācijas darba vadītāju
- Implementācijas rīka[3] studēšana
- Implementācijas rīkā izmantotās literatūras[7] studēšana

7. DARBIETILPĪBAS NOVĒRTĒJUMS

Darbietilpības rezultāts

Darba ietilpības aprēķins veikts izmantojot *Azure DevOps* un *Microsoft Excel*. Zemāk redzama tabula, kas satur datus par paveiktā darba izpildei patērēto darba dienu skaitu, pieņemot, ka viena darba diena ir astoņas stundas.

7.1 tabula

Patērēto darba dienu skaits izstrādes laikā

Paveiktais darbs	Dienu skaits
Dokumentācijas izstrāde	
Konsultācija ar pasūtītāju	10
Prasību specificēšana	10
Prasību specifikācijas saskaņošana	3
Projektējuma izveide	10
Konsultācijas ar kolēģiem	3
Projekta izstrādes plānošana	1
<i>Coretopic</i> izpēte	5
Testu dokumentēšana	3
Funkcionālo klašu programmēšana, vienībtestēšana	
Trainer	0.5
Factory	0.5
CorExConfiguration	5
CorExExtractor	20
Automātisko vienībutestu izstrāde, testēšana	
TestCorExExtractor	3
TestCorExConfiguration	2.5
TestFactory	0.5
KOPĀ	77

- Kopā patērētās dienas: 77
- Kopā patērētās darba stundas: 616
- Stundu skaits vienā personmēnesī: 152
- *CorEx*[7] implementācija veikta aptuveni: 4 personmēnešos

COCOMO darbietilpības mērījums

Izmantojot COCOMO[4][5] darbietilpības mērīšanas metodi tika aprēķināts, ka EAF[4] vienāds ar 1.15, un tā kā ir trīs vienkārši ievadi un trīs vienkārši izvadi, tad funkcijpunktu (FP) skaits ir vienāds ar 24[5]. Tā kā projekts tiek veidots izmantojot Python valodu, tad LOC/FP (LOC - *lines of code* jeb koda rindiņas) koeficients ir 67. Lai aprēķinātu rindiņu skaitu, sareizina FP ar LOC/FP, kas ir $24 \times 67 = 1,608$. Tātad KLOC (kiloLOC jeb tūkstots koda rindiņu) ir vienāds ar 1,6. Beigu beigās aprēķinot personmēnešu skaitu izmanto formulu $a \cdot (\text{KLOC}^b) \cdot \text{EAF}$ [4], kur a vienāds ar 3.0 un b ir 1.12, iegūst, ka personmēnešu skaits ir **5.84** personmēneši.

NOBEIGUMS

CorEx[7] implementācija ir sekmīgi sekojot prasību specifikācijai un ievērojot uzstādīto terminu. Ir gan palicis izstrādātā projekta lietošana un lietotāju testēšana, kā arī ir iespējams veikt dažus uzlabojumus un papildinājumus. Šādi papioldinājumi varētu būt *Corextopic*[3] bibliotēkā piedāvātās anchored trenēšanas implementācija, kas nodrošinātu rezultātu precīzāku iegūšanu sekojot lietotāja uzstādījumiem un papildus ievaddatiem.

Izstrādes laikā tika apgūts kā projekta pārvaldības un projekta dokumentācijas pamati. Grūtākais bija izprast, kā pasūtījam jau eksistējošā sistēma (*MMTA*) darbojās un kā izpildīt pasūtītāja prasības.

Tika iegūta nozīmīgi svarīga pieredze programmēšanas valodas *Python 3.5.5* pielietojumā, kā arī mašīnamācīšanās izstrādes pamatos, konkrētāk teksta analīzes sfērā. Kā arī tika iegūta pieredze veidojot vienībtestus izmantojot *Python 3.5.5* satvaru *unittest*.

Šobrīd turpinās *CorEx*[7] implementācijas uzlabojumi un funkcionāli papildinājumi, kā arī tiks veikti papildinājumi un citu funkcionalitāšu implementēšana *MMTA* sistēmā.

IZMANTOTĀ LITERATŪRA

Grāmatas

1. Pressman R. *Software engineering: a practitioner's approach*. New York: McGraw-Hill Higher Education, 2010.

Tīmekļa resursi

2. Azure DevOps [tiešsaiste] [atsauce 23.05.2019] Pieejams internetā: <https://dev.azure.com>
3. Corextopic [tiešsaiste] [atsauce 23.05.2019] Pieejams internetā: https://github.com/gregversteeg/corex_topic
4. Cocomo [tiešsaiste] [atsauce 23.05.2019] Pieejams internetā: <https://en.wikipedia.org/wiki/COCOMO>
5. Cocomo lietošanas piemērs [tiešsaiste] [atsauce 23.05.2019] Pieejams internetā: <https://estudijas.lu.lv/mod/resource/view.php?id=128556>
6. What is Scrum? [tiešsaiste] [atsauce 23.05.2019] Pieejams internetā: <https://www.scrum.org/resources/what-is-scrum>
7. Gallagher, Ryan J., Kyle Reing, David Kale, and Greg Ver Steeg. "Anchored Correlation Explanation: Topic Modeling with Minimal Domain Knowledge." Transactions of the Association for Computational Linguistics (TAACL), 2017. pp.529-531 [tiešsaiste] [atsauce 23.05.2019] Pieejams internetā: <https://www.transacl.org/ojs/index.php/tacl/article/view/1244>,
8. Draw.io [tiešsaiste] [atsauce 23.05.2019] Pieejams internetā: <https://www.draw.io/>
9. CountVectorizer dokumentācija [tiešsaiste] [atsauce 23.05.2019] Internetā: https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html
10. Letonika [tiešsaiste] [atsauce 23.05.2019] Pieejams internetā: <https://www.letonika.lv/groups/default.aspx?r=1107&q=v%C4%81rdn%C4%ABca&id=2071241&g=1>

Latvijas Valsts standarti

11. Latvijas Valsts Standarts LVS 68:1996 "Programmatūras prasību specifikācijas ceļvedis",
12. Latvijas Valsts Standarts LVS 72:1996 "Ieteicamā prakse programmatūras projektējuma aprakstīšanai"

Kvalifikācijas darbs „*Papildu funkciju izstrāde mašīnmācīšanās teksta analizatoram*” izstrādāts Latvijas Universitātes Datorikas fakultātē.

Ar savu parakstu apliecinu, ka darbs izstrādāts patstāvīgi, izmantoti tikai tajā norādītie informācijas avoti un iesniegtā darba elektroniskā kopija atbilst izdrukai.

Autors: *Mikus Kalniņš* _____ .05.2019.

Rekomendēju darbu aizstāvēšanai

Darba vadītājs/a: *prof. Jānis Zuters* _____ .05.2019.

Recenzents: *Jānis Vempers, ZZ Dats SIA, projektu pārvaldnieks*

Darbs iesniegts 28.05.2018.

Kvalifikācijas darbu pārbaudījumu komisijas sekretārs: *Darja Solodovņikova* _____

Darbs aizstāvēts kvalifikācijas darbu pārbaudījuma komisijas sēdē

____.06.2019. prot. Nr. _____

Komisijas sekretārs(-e): _____