

LATVIJAS UNIVERSITĀTE  
DATORIKAS FAKULTĀTE

**AUTOMĀTISKA E-PASTU APSTRĀDES UN ATBILDES  
SNIEGŠANAS SISTĒMA**

KVALIFIKĀCIJAS DARBS

Autors: **Edvards Bukovskis**

Studenta apliecības Nr.: eb18081

Darba vadītājs: Dr. dat. Jānis Zuters

RĪGA 2021

## Anotācija

Kvalifikācijas darbā tika izstrādāts Automātiskās E-pastu Atbildes Sniegšanas Sistēmas (turpmāk - AEAASS) projektējums, programmatūras specifikācija un daļēja sistēmas programmatūras realizācija. Sistēma un tās koncepti tika veidoti reāla uzņēmuma projekta ietvaros, un šis programmatūras specifikācijas dokuments ir pamats turpmākai sistēmas programmatūras attīstībai un uzlabošanai nākotnē.

AEAASS ir risinājums reālai darbinieku darba slodzes un šī brīža COVID-19 pandēmijas izraisīto seku problēmām, uzņēmuma interesēs. Sistēmas būtība ir lielo datu izmantošana konkrētu darbību automatizācijai, izmantojot tādus elementus kā: dabiskā valodas apstrāde (NLP) un mašīnāpmācība (ML).

Atslēgvārdi: automatizācija, datu apstrāde, lielle dati, mašīnāpmācība, dabiskā valodas apstrāde, Python

# Automatic E-mail Processing and Response System

## Abstract

In the qualification work, the design of the Automatic E-mail Response System (hereinafter - AEAASS), software specification and partial implementation of the system software were developed. The system and its concepts were developed within the framework of a real company project, and this software specification document is the basis for further development and improvement of the system software in the future.

AEAASS is a solution to the real problems of employee workload and the consequences of the current COVID-19 pandemic, in the interests of the company. The essence of the system is the use of big data to automate specific activities using elements such as: natural language processing (NLP) and machine learning (ML).

Keywords: Automation, Data Processing, Big Data, Machine Learning, Natural Language Processing, Python

# Saturs

1	IEVADS	7
1.1	Nolūks	7
1.2	Darbības sfēra	7
1.3	Saistība ar citiem dokumentiem:	7
1.4	Pārskats:	7
2	Vispārējs apraksts	8
2.1	Esošā stāvokļa apraksts	8
2.2	Pasūtītājs	8
2.3	Produkta perspektīva	8
2.4	Sistēmas funkcionalitāte	8
2.5	Sistēmas lietotāji	8
2.6	Vispārējie ierobežojumi	9
2.7	Pieņemumi un atkarības	9
3	Programmatūras prasību specifikācija	10
3.1	Datu plūsma	10
3.2	Funkcionālās prasības	11
3.2.1	Moduļu sadalījums un apraksts	11
3.2.2	Teksta apstrādes moduļa apraksts	11
3.2.3	Teksta satura identificēšanas modulis	14
3.3	Nefunkcionālās prasības	16
3.3.1	Uzturamība	16
3.3.2	Lietojamība	16
3.3.3	Veiktspēja	16
3.3.4	Izvietojamība	16
3.3.5	Drošība	16
4	Programmatūras projektējuma apraksts	17
4.1	Ievads	17
4.1.1	Nolūks	17

4.1.2	Darbības sfēra	17
4.1.3	Saistībā ar citiem dokumentiem	17
4.2	Moduļu dekompozīcija	17
4.3	ML modeļa dekompozīcija	19
4.3.1	Datu attīrīšanas soļu apraksts	20
4.3.2	Datu attīrīšanas soļu reāls vizualizējums ar piemēru <b>pozitīviem</b> datiem	25
4.3.3	Datu attīrīšanas soļu reāls vizualizējums ar piemēru <b>negatīviem</b> datiem.	26
4.4	Daļējs funkciju projektējums	28
5	Testēšana	29
6	Projekta organizācija	32
7	Konfigurācijas pārvaldība	33
8	Darbietilpības novērtējums	34
9	Secinājumi un autora komentāri	35
10	Izmantotā literatūra un avoti	36
	Pielikumi	37

## Vārdnīca

<b>NLP</b>	Dabiskā valodas apstrāde (no angļu <i>Natural Language Processing</i> )
<b>PPS</b>	Programmatūras prasību specifikācija
<b>PPA</b>	Programmatūras projektējuma apraksts
<b>AEAASS</b>	Automātiska e-pastu apstrādes un atbildes sniegšanas sistēma
<b>ML</b>	Mašīnāpmācība (no angļu <i>Machine Learning</i> )
<b>Python</b>	Augsta līmeņa interpretējama objektorientētā programmēšanas valoda
<b>SMTP</b>	Vienkāršais pasta pārsūtīšanas protokols
<b>Firewall</b>	Ugunsmūris

# **1 IEVADS**

## **1.1 Nolūks**

Dokuments ir sagatavots kvalifikācijas darba un uzņēmuma projekta ietvaros. Dokumentā tiek aprakstīta AEAASS struktūra un funkcionalitāte, kā arī, dokumentns sevī ietver pasūtītāju gatavā produkta redzējumu un tā daļēju tehnisko izpildījumu, kas turpmāk uzņēmumam būs pamats turpmākai sistēmas attīstībai.

Darbs tika veikts saistībā ar iepriekšējo praksi uzņēmumā SIA “Tet”, lai risinātu problēmu uzņēmuma interesēs.

## **1.2 Darbības sfēra**

AEAASS ir paredzēts intergrēt uzņēmuma klientu apkalpošanas un palīdzības nodaļā, lai atvieglotu darba slodzi un virzītos uz kopēju sistēmu automatizēšanu. AEAASS būtu apakšsistēma jau esošai uzņēmuma sistēmai.

## **1.3 Saistība ar citiem dokumentiem:**

Dokumenta PPS noformēšanā ievērotas standarta LVS 68:1996 prasības.

AEAASS PPS ir cieši saistīts ar šīs pašas sistēmas PPA.

Dokumenta PPS un PPA kalpo kā atsauce turpmākiem sistēmas uzlabojumiem un citām šīs sistēmas novirzēm uzņēmumā.

## **1.4 Pārskats:**

Pirmajā daļā ir ievadinformācija un vispārējs apraksts, kas satur dokumenta nolūku, definīciju skaidrojumu, Saistību ar citiem dokumentiem un sistēmas darbības sfēru.

Otrā daļa sniedz pārskatu par sistēmas esošo stāvokli, pasūtītāju, sistēmas perspektīvu un vispārēju sistēmas funkcionalitāti un būtību.

Trešajā daļā tiek aprakstītas funkcionālās prasības un nefunkcionālās prasības

Ceturtajā daļā PPA tiek sīkāk aprakstīta nozīmīgu sistēmas struktūrvienību uzbūve un darbošanās princips.

Piektajā daļā tiek aprakstīta funkciju testēšana.

## **2 VISPĀRĒJS APRAKSTS**

### **2.1 Esošā stāvokļa apraksts**

Ņemot vērā uzņēmuma pieaugošo klientu skaitu, kā arī, pašreizējo situāciju pandēmijas laikā, klientu apkalpošanas un palīdzības nodaļai rodas lielāka noslodze uz esošo darbinieku skaitu. Lai risinātu šo radušos problēmu, rodas nepieciešamība automatizēt atsevišķas, monotonas darbības. Risinājumam ir jābūt pēc iespējas automatizētam, lai mazinātu darbinieku noslogojumu un veicinātu augstāku produktivitāti, tāpēc tiek veidots projekts, kura ietvaros tiek izstrādāta apakšsistema, kas papildinās esošo sistemu.

### **2.2 Pasūtītājs**

Faktiskais sistēmas pasūtītājs ir klientu apkalpošanas un palīdzības nodaļa uzņēmumā.

### **2.3 Produkta perspektīva**

Sistēma pati par sevi ir neatkarīga, tāpēc to ir iespējams pielāgot un modificēt citām vajadzībām. Konkrētajā gadījumā, kad sistēma kalpo kā apakšsistēma jau esošai sistēmai, tai ir iespējams pievienot papildu funkcionalitātes.

### **2.4 Sistēmas funkcionalitāte**

1. Teksta apstrāde
2. Teksta satura atpazīšana
3. Atbildes sniegšana, balstoties uz teksta saturu
4. ML modeļa trenēšana

### **2.5 Sistēmas lietotāji**

Sistēmas netiešie lietotāji ir uzņēmuma klientu apkalpošanas un palīdzības nodaļas darbinieki, jo sistēma darbojas automātiski un tiešas mijiedarbības ar fiziskām personām nav paredzētas.



## **2.6 Vispārējie ierobežojumi**

Sistēmu ir nepieciešams integrēt SMTP serverī, kā arī, nepieciešams interneta savienojums.

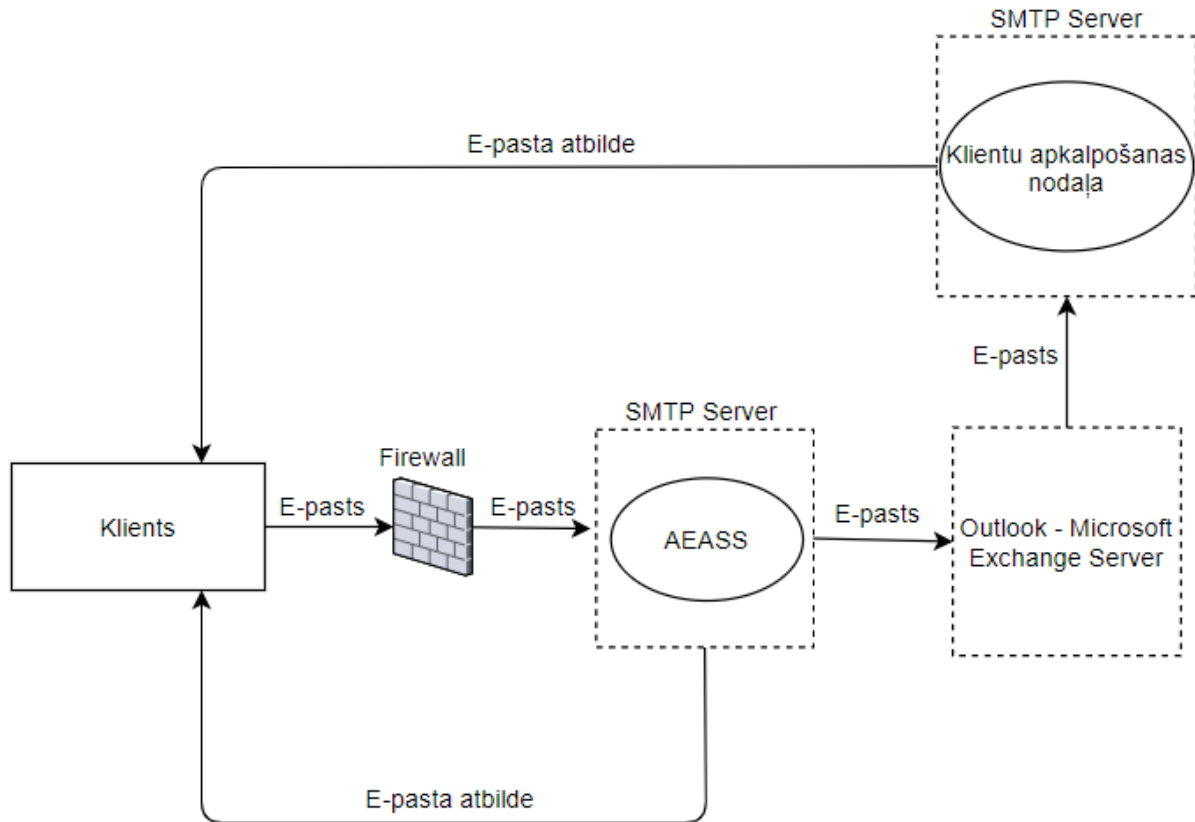
## **2.7 Pieņēmumi un atkarības**

AEAASS izstrādātie pieņēmumi, kuri nepieciešami sistēmas darbībai:

- Sistēma ir citas sistēmas apakšsistēma

### 3 PROGRAMMATŪRAS PRASĪBU SPECIFIKĀCIJA

#### 3.1 Datu plūsma

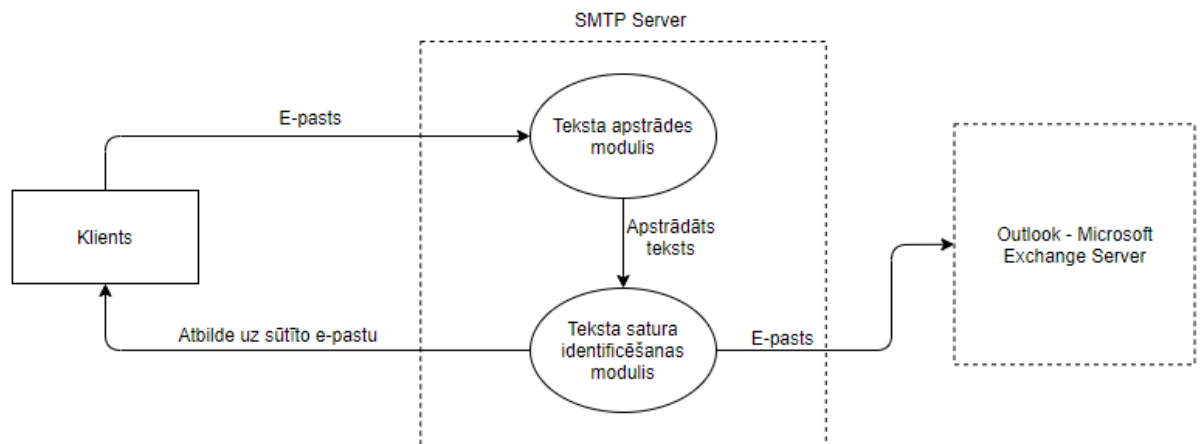


Attēls 3.1 0. līmeņa datu plūsmas diagramma

AEASS ir kā starpposms (skat. Attēls 3.1), kas apstrādā klientu e-pastus pirms tie nonāk klientu apkalpošanas un palīdzības nodaļas sistēmā. AEASS tiek integrēts SMTP serverī, kur ar apmācītu ML modeli tiek veikta e-pastu filtrācija. Tālāk tiek realizēts “luksofora” princips – ja e-pasts atbilst AEASS kritērijiem, tad tiek automātiski nosūtīta atbilde atpakaļ klientam, pretējā gadījumā, e-pasts tiek nodots tālāk uz manuālu apstrādi klientu apkalpošanas un palīdzības nodaļai.

## 3.2 Funkcionālās prasības

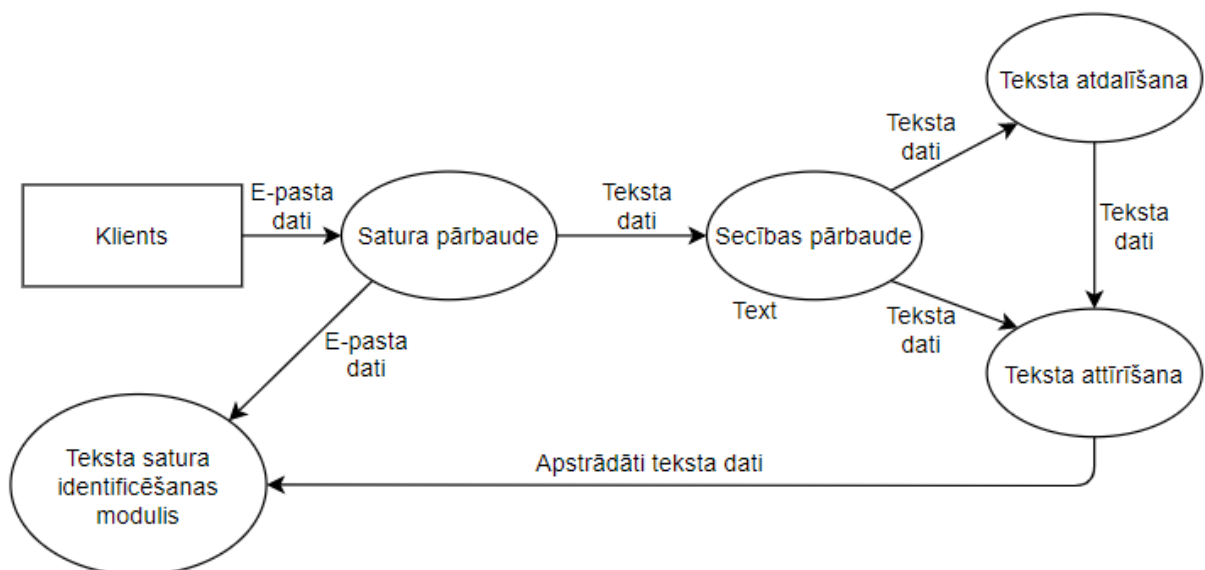
### 3.2.1 Moduļu sadalījums un apraksts



Attēls 3.2 1. līmeņa datu plūsmas diagramma

AEASS satur divus moduļus (skat. Attēls 3.2) – Teksta apstrādes moduli un teksta satura pārbaudes moduli. Teksta apstrādes modulis atbild par ienākošā e-pasta teksta formāta sagatavošanu padošanai ML modelis, kas ir iekļauts teksta satura pārbaudes modulī. Teksta satura identificēšanas modulis atbild par teksta satura būtības pārbaudi, atbilžu sniegšanu, vai e-pasta nosūtīšanu tālākai apstrādei uz klientu apkalpošanas centra sistēmu.

### 3.2.2 Teksta apstrādes moduļa apraksts



Attēls 3.3 Teksta apstrādes moduļa datu plūsmas diagramma

Teksta apstrādes modulis (skat. Attēls 3.3) sastāv no četrām funkcijām – Satura pārbaudes, e-pasta secības pārbaudes, teksta atdalīšanas un teksta attīrīšanas funkcijām. Katras funkcijas mērķis šajā modulī ir apstrādāt teksta datus tādā veidā, lai Teksta satura identificēšanas modulis tos spētu atpazīt un pielietot ML modelim.

Tabula 3.1

**Satura pārbaudes funkcijas apraksts**

<b>Mērķis:</b>
Atlasīt gadījumus, kuros e-pasta saturā tiek sūtīts pielikums, piemēram, dokuments vai attēls
<b>Ievaddati:</b>
SMTP serverī saņemtais e-pasts
<b>Apstrāde:</b>
Tiek pārbaudīts vai e-pasts satur tikai simbolu virkni, vai arī pielikumu/-us
<b>Izvaddati:</b>
<p>Gadījumā, ja e-pasts nesatur pielikumus:</p> <ul style="list-style-type: none"> <li>• Dati tiek nodoti tālāk secības pārbaudes funkcijai</li> </ul> <p>Gadījumā, ja e-pasts satur pielikumu:</p> <ul style="list-style-type: none"> <li>• Tam tiek piešķirta speciāla iezīme un e-pasts tiek nodots Teksta satura identificēšanas modulim</li> </ul>

Tabula 3.2

## Secības pārbaudes funkcija

<b>Mērķis:</b>
Atšķirt e-pastus, kurus klients raksta pirmo reizi un e-pastus, kas ir daļa no pastāvošas sarakstes.
<b>Ievaddati:</b>
Teksta dati (simbolu virkne), kurus atgriež funkcija content_check
<b>Apstrāde:</b>
Tiek pārbaudīts vai simbolu virknē ir atrodamas iezīmes, kas liecinātu par to, ka e-pasta teksts nav pirmreizējs, t.i. vai e-pasts nav sastāvdaļa no esošas sarakstes ar klientu apkalpošanas un palīdzības nodaļu.
<b>Izvaddati:</b>
Teksta attīrīšanas vai teksta atdalīšanas funkciju izsaukums

Tabula 3.3

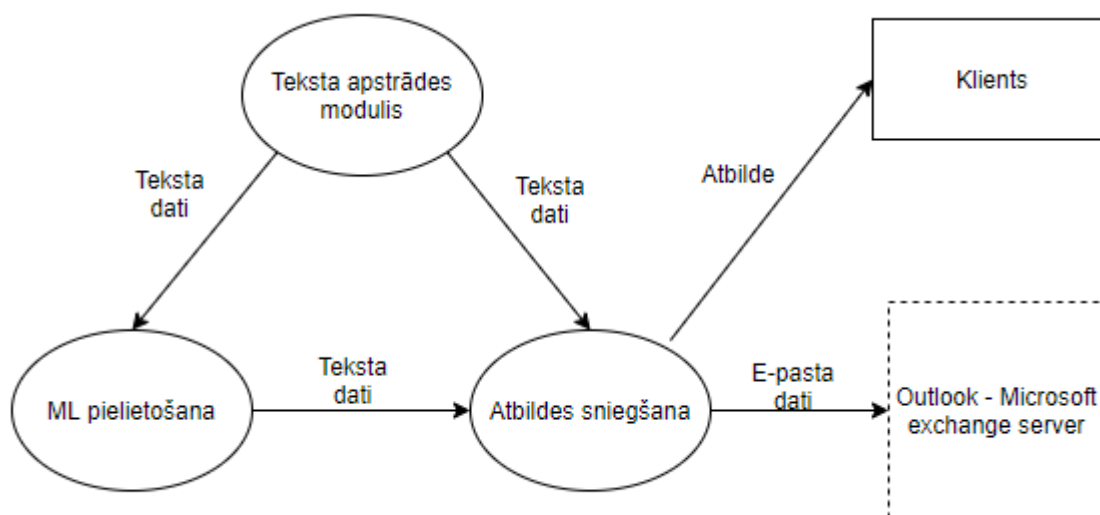
## Teksta atdalīšanas funkcijas apraksts

<b>Mērķis:</b>
Atdalīt nevajadzīgās simbolu virknes daļas, gadījumā, ja e-pasts nav pirmreizējs, t.i. ja e-pasts ir sastāvdaļa no sarakstes
<b>Ievaddati:</b>
Teksta dati (simbolu virkne), kurus atgriež funkcija content_check
<b>Apstrāde:</b>
E-pastu sarakstes struktūrai ir konkrētas, statiskas iezīmes, pēc šīm iezīmēm tiek atdalīts nepieciešamais teksts.
<b>Izvaddati:</b>
Simbolu virkne, kurā ir konkrēts klienta sūtīts e-pasta jautājums.

## Teksta attīrīšanas funkcijas apraksts

<b>Mērķis:</b>
Apstrādāt simbolu virkni, lai to spētu atpazīt un identificēt ML modelis Teksta satura identificēšanas modulī
<b>Ievaddati:</b>
Simbolu virkne no funkcijas split_important vai pirmreizējā simbolu virkne.
<b>Apstrāde:</b>
<ol style="list-style-type: none"> <li>1. Simbolu virkne tiek attīrīta no nesvarīgas informācijas, kura nemaina teksta galveno domu vai būtību, piemēram, simbolu virkne tiek attīrīta no e-pastu adresēm, biežākajiem stopvārdiem.</li> <li>2. Simbolu virkne tiek transformēta mazajos burtos un tiek noņemtas garumzīmes.</li> <li>3. Simbolu virknes, kuras nesatur latīņu alfabēta burtus, tiek speciāli iezīmētas.</li> </ol>
<b>Izvaddati:</b>
Apstrādāta un attīrīta simbolu virkne

## 3.2.3 Teksta satura identificēšanas modulis



Attēls 3.4 Teksta identificēšanas moduļa datu plūsmas diagramma

Teksta identificēšanas modulis (skat. Attēls 3.4) sastāv no divām pamatfunkcijām – ML pielietošanas un atbildes sniegšanas. ML pielietošanas funkcija atbild par teksta klasificēšanu, izmantojot apmācītu ML modeli. Pēc attiecīgā klasifikatora, tiek sniegta automātiska e-pasta atbilde klientam vai arī e-pasts tālāk tiek nodots uz klientu apkalpošanas un palīdzības centra sistēmu.

Tabula 3.5

#### ML pielietošanas funkcijas apraksts

<b>Mērķis:</b>
Identificēt un iezīmēt tekstu pēc satura būtības
<b>Ievaddati:</b>
Apstrādāta simbolu virkne no teksta apstrādes moduļa
<b>Apstrāde:</b>
Tiek pielietots apmācīts ML modelis, kas identificē un klasificē simbolu virkni
<b>Izvaddati:</b>
Klasificēta simbolu virkne

Tabula 3.6

#### Atbildes sniegšanas funkcijas apraksts

<b>Mērķis:</b>
Sniegt iepriekš sagatavotu atbildi klientam vai nodot e-pastu tālākai apstrādei uzņēmuma klientu apkalpošanas un palīdzības nodaļas sistēmai.
<b>Ievaddati:</b>
E-pasta dati, klasificēta simbolu virkne.
<b>Apstrāde:</b>
Funkcijai ir noteiktas darbības pie konkrētām iezīmēm, un attiecībā no iezīmes tiek sniegta vai nu atbilde klientam, vai e-pasts tiek nodots tālākai apstrādei.
<b>Izvaddati:</b>
E-pasta atbilde

### **3.3 Nefunkcionālās prasības**

#### **3.3.1 Uzturamība**

Programmas kodam ir jābūt labi strukturētam, komentētam un viegli uztveramam, jāievēro labs programmēšanas stils

#### **3.3.2 Lietojamība**

Sistēma nav paredzēta tiešai mijiedarbībai ar klientu apkalpošanas un palīdzības nodaļas darbiniekiem, sistēma tiek integrēta jau esošā e-pastu apmaiņas sistēmā kā starpposms. Ja rodas nepieciešamība pie sistēmas koda rediģēšanas vai sistēmas koda modificēšanas citiem nolūkiem uzņēmumā, tad attiecīgi konkrētām personām ir jānodrošina pieeja šādai iespējai.

#### **3.3.3 Veiktspēja**

Sistēmai ir jādarbojas ar augstu precizitāti, pretējā gadījumā paredzētais darba slogs netiks atvieglināts, bet papildu noslogots.

#### **3.3.4 Izvietojamība**

Tiek paredzēts, ka sistēmu ir viegli modificēt, un sistēmu vai tās daļas ir iespējams viegli pielāgot citiem uzņēmuma nolūkiem.

#### **3.3.5 Drošība**

Sistēmas datiem ir jābūt droši glabātiem, lai tiem nav iespējams piekļūt no ārējām sistēmām. Piekļuvi sistēmas datiem, pēc nepieciešamības, drīkst sniegt tikai kvalificētam personālam konkrētiem un skaidriem nolūkiem.



## **4 PROGRAMMATŪRAS PROJEKTĒJUMA APRAKSTS**

### **4.1 Ievads**

#### **4.1.1 Nolūks**

AEAASS programmatūras projektējuma apraksts sniedz padziļinātu informāciju par sistēmas elementu uzbūvi, struktūru un darbības principiem. PPA tiek aprakstīta ML modeļa uzbūve, modeļa trenēšana ar datiem un darbības princips, kas sevī ietver lielo datu apstrādi ar NLP metodēm.

#### **4.1.2 Darbības sfēra**

AEAASS PPA ir paredzēts padziļinātai informācijai par sistēmas struktūrvienībām. Šis dokuments var tikt izmantots kā pamats turpmākai AEASS attīstībai vai citām tās novirzēm, kas iekļauj AEAASS elementus.

#### **4.1.3 Saistībā ar citiem dokumentiem**

PPA ir cieši saistīts ar AEAASS PPS aprakstu. PPA ir izstrādāts pēc standarta "LVS 72:1996 Ieteicamā prakse programatūras projektējuma aprakstīšanai".[3]

### **4.2 Moduļu dekompozīcija**

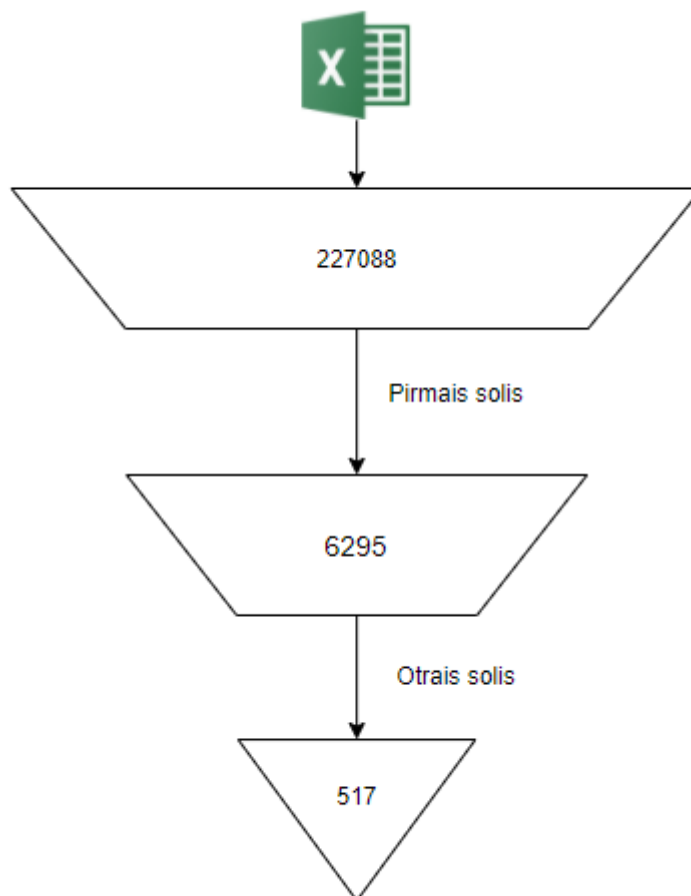
Dekompozīcijā tiek sīkāk aprakstīti sistēmā esošo moduļu funkcijas un elementi. Tabulā tiek parādītas konkrētas sistēmas funkcijas, to apraksts un iekš kura moduļa funkcija tiek izmantota. (skat. Tabula 4.1).

Tabula 4.1

<b>Funkcija</b>	<b>Apraksts</b>	<b>Teksta apstrādes modulis</b>	<b>Teksta satura identificēšanas modulis</b>
Content_check(string)	<i>Atgriež vērtības True vai False , Pārbauda vai e-pastā ir iekļauts pielikums.</i>	+	
sequence_check(string)	<i>Atgriež vērtības True vai False, pārbauda vai e-pasts ir pirmreizējs.</i>	+	
split_important(string)	<i>Atgriež simbolu virkni, apstrādā simbolu virkni, atstājot nepieciešamo informāciju.</i>	+	
text_clean(string)	<i>Atgriež simbolu virkni, apstrādā un attīra tekstu, sagatavojot to ML modeļa atpazīstamajā formā.</i>	+	
modelCompile(string)	<i>Atgriež "vārdnīcu", identificē un klasificē attīrītu simbolu virkni.</i>		+
send_response(string[])	<i>Sūta sagatavotu e-pasta atbildi klientam vai nodod e-pastu tālākai apstrādei</i>		+

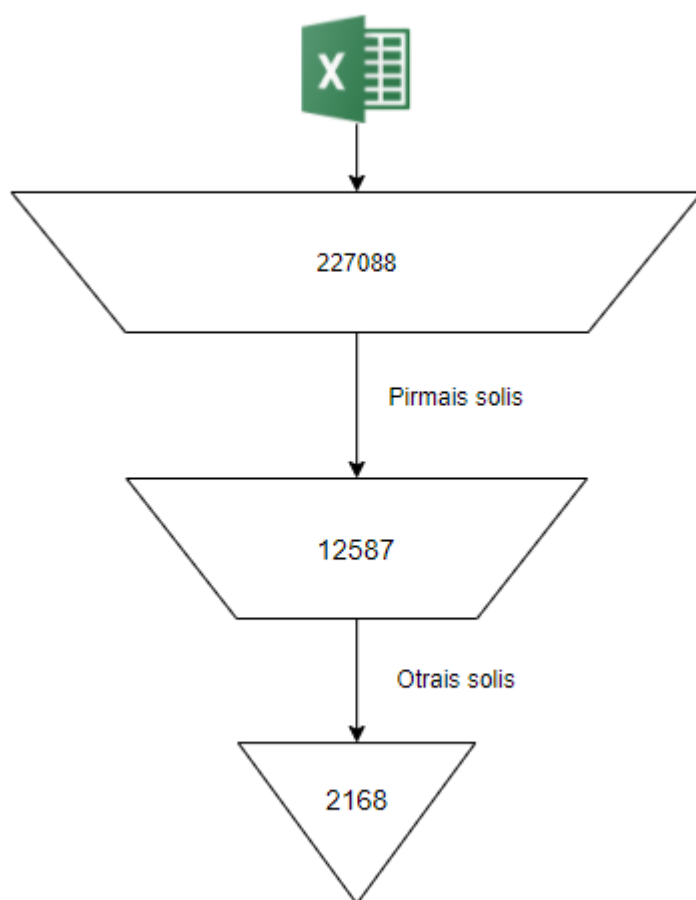
### 4.3 ML modeļa dekompozīcija

Dekompozīcijā tiek aprakstīti un parādīti veiktie soļi ML modeļa trenēšanā [5]. Dekompozīcijā tiek arī sīkāk aprakstīta teksta apstrādes loģika. Dotajos piemēros tiek aprakstīta klientu identificēšanās iezīmes izveide. Tiek sagatavoti pozitīvie dati (skat. Attēls 4.1) un negatīvie dati (skat. Attēls 4.2). Dati sākotnēji tiek ielasīti Microsoft Excel dokumentā no uzņēmuma datu noliktavas. Datu būtība ir izsūtītie un saņemtie e-pasti divu gadu garumā.



*Attēls 4.1. Pozitīvo datu filtrācijas diagramma*

(skat. Attēls 4.1.) Attēlā tiek parādīts pieejamo datu skaits un tā attiecība pret datu attīrīšanas veiktajiem soļiem pozitīvo datu atlasīšanā.



*Attēls 4.2. Negatīvo datu filtrācijas diagramma*

(skat. Attēls 4.2.) Attēlā tiek parādīts pieejamo datu skaits un tā attiecība pret datu attīrīšanas veiktajiem soļiem negatīvo datu atlasīšanā.

#### 4.3.1 Datu attīrīšanas soļu apraksts

Katras diagrammās (skat. Attēls 4.1, 4.2) Solis sevī iekļauj simbolu virknes attīrīšanas elementus, izmantojot NLP metodes. (sīkāk skat. Tabulas 4.1, 4.2)

Datu attīrīšanas soļu apraksts pozitīviem datiem

Solis	Darbības
<i>Pirmais solis</i>	<p>Būtība:</p> <p>Tiek atlasīti klientu jautājumi uz kuriem, lai sniegtu atbildi, <b>IR</b> nepieciešams identificēties.</p> <p>1. Tiek atlasīti tikai tādi dati (šajā gadījumā e-pasti), kuros <b>IR</b> minētas konkrētas frāzes, piemēram:</p> <ul style="list-style-type: none"> <li>• str.contains('Klienta numurs')</li> <li>• str.contains('Personas kods')</li> </ul> <p>Būtiski, šajā solī netiek izmantoti regeksi, lai paātrinātu darbību.</p> <p>2. Tālāk, tiek atlasīti tikai tādi dati, kuri nāk no privātajiem klientiem, respektīvi, tiek meklētas privātas e-pastu adreses, piemēram:</p> <ul style="list-style-type: none"> <li>• Gmail.com</li> <li>• Inbox.lv</li> <li>• Yahoo.com</li> </ul> <p>3. Tiek apstrādāti ne-pirmreizēji e-pasti:</p> <ul style="list-style-type: none"> <li>• Tiek meklētas sakarības starp e-pasta sarakstes elementiem, konkrētajā gadījumā, tiek meklēta simbolu virkne "--Original Message--", kas atdala klienta pirmreizējo e-pasta tekstu.(skat. tālāk attēlu piemēros 4.3, 4.4, 4.5)</li> <li>• Nepieciešamās teksta daļas tiek atdalītas, šajā gadījumā: <ul style="list-style-type: none"> <li>○ str.split("--Original Message--")</li> </ul> </li> </ul>

Solis	Darbības
<i>Otrais solis</i>	<p>Būtība:</p> <p>Dati tiek apstrādāti pēc NLP[4] metodēm, lai tos izmantotu ML modeļa trenēšanai.</p> <ol style="list-style-type: none"> <li>1. Simbolu virknes tiek pārveidotas mazajos burtos</li> <li>2. Tiek attīrītas tiešsaistes adreses un e-pastu adreses</li> <li>3. Tiek attīrīti pieklājības vārdi/frāzes un sasveicināšanās vārdi/frāzes (Paldies, labrīt, uz redzēšanos, ar cieņu, u.tml.)</li> <li>4. Garumzīmes un mīkstinājuma zīmes tiek nomainītas uz attiecīgo burtu bez garumzīmes vai mīkstinājuma zīmes.</li> <li>5. Tiek attīrīti jebkādi simboli izņemot latīņu burtus (kirilicas simboli, @, (, /, &lt;, u.tml.)</li> <li>6. Tiek attīrītas nestandarta frāzes, kas tiek automātiski izveidotas atkarībā no ierīces vai e-pastu sistēmas. ('sent from', 'message created on', u.tml.)</li> <li>7. Tiek attīrīti stopvārdi (ka ,lai, jo bet, u.tml.)[2]</li> <li>8. Tiek attīrīti "vientuļie" simboli.</li> </ol>

Datu attīrīšanas soļu apraksts negatīviem datiem

Solis	Darbības
<i>Pirmais solis</i>	<p>Būtība:</p> <p>Tiek atlasīti klientu jautājumi uz kuriem, lai sniegtu atbildi, <b>NAV</b> nepieciešams identificēties.</p> <p>1. Tiek atlasīti tikai tādi dati (šajā gadījumā e-pasti), kuros <b>NAV</b> minētas konkrētas frāzes, piemēram:</p> <ul style="list-style-type: none"> <li>• str.contains('Klienta numurs')</li> <li>• str.contains('Personas kods')</li> </ul> <p>Būtiski, šajā solī netiek izmantoti regeksi, lai paātrinātu darbību.</p> <p>2. Tālāk, tiek atlasīti tikai tādi dati, kuri nāk no privātajiem klientiem, respektīvi, tiek meklētas privātas e-pastu adreses, piemēram:</p> <ul style="list-style-type: none"> <li>• Gmail.com</li> <li>• Inbox.lv</li> <li>• Yahoo.com</li> </ul> <p>3. Tiek apstrādāti ne-pirmreizēji e-pasti:</p> <ul style="list-style-type: none"> <li>• Tiek meklētas sakarības starp e-pasta sarakstes elementiem, konkrētajā gadījumā, tiek meklēta simbolu virkne "--Original Message--", kas atdala klienta pirmreizējo e-pasta tekstu.(skat. tālāk attēlu piemēros 4.3, 4.4, 4.5)</li> <li>• Nepieciešamās teksta daļas tiek atdalītas, šajā gadījumā: <ul style="list-style-type: none"> <li>○ str.split("--Original Message--")</li> </ul> </li> </ul>

Solis	Darbības
<i>Otrais solis</i>	<p>Būtība:</p> <p>Dati tiek apstrādāti pēc NLP metodēm[4], lai tos izmantotu ML modeļa trenēšanai.</p> <ol style="list-style-type: none"> <li>1. Simbolu virknes tiek pārveidotas mazajos burtos</li> <li>2. Tiek attīrītas tiešsaistes adreses un e-pastu adreses</li> <li>3. Tiek attīrīti sasveicināšanās vārdi/frāzes (Paldies, labrīt, uz redzēšanos, ar cieņu, u.tml.)</li> <li>4. Garumzīmes un mīkstinājuma zīmes tiek nomainītas uz attiecīgo burtu bez garumzīmes vai mīkstinājuma zīmes.</li> <li>5. Tiek attīrīti jebkādi simboli izņemot latīņu burtus (kirilicas simboli, @, (, /, &lt;, u.tml.)</li> <li>6. Tiek attīrītas nestandarta frāzes, kas tiek automātiski izveidotas atkarībā no ierīces vai e-pastu sistēmas. ('sent from', 'message created on', u.tml.)</li> <li>7. Tiek attīrīti stopvārdi (ka ,lai, jo bet, u.tml.)[2]</li> <li>8. Tiek attīrīti "vientuļie" simboli.</li> </ol>



#### 4.3.2 Datu attīrīšanas soļu reāls vizualizējums ar piemēru **pozitīviem** datiem

*Jūtīga informācija tiek aizstāta ar simbolu 'x'.*

Labdien! Lūdzu atsūtiet linku paroles atjaunošanai. On 19.10.2019 18:02, Tet wrote: >> Sveiki! >> Paldies par iesūtīto vēstuli! >> Nodevām atbildīgajiem speciālistiem, pārbaudei. >> Atbildes uz jautājumiem aicinām meklēt arī mūsu mājas lapas Palīdzības > sadaļā - [http://palidziba.tet.lv/lat/palidziba\\_majai/](http://palidziba.tet.lv/lat/palidziba_majai/). > Ērtākai ikdienai aicinām ar mums sazināties čatā. >> Tavš Tet >>>> --Original Message-- > From: xxx@gmail.com > Date: 19.10.2019 13:31 > To: tet@tet.lv > Subject: Re: atbloķēt e-pastu [#234449] >> xxx xxx > xxxxx-xxxxx> Viestura laukums, valmiera | On 19 Oct 2019, 13:26 +0300, Tet <tet@tet.lv>, wrote: >> Sveiki! >> Paldies par Jūsu vēstuli! >> Lai pārbaudīt radušos situāciju ar Jūsu e-pastu, lūgums, saistībā ar > personas datu aizsardzību, norādīt: > 1. klienta numuru vai personas kodu > 2. vārds, uzvārds > 3. pakalpojuma adresi. >> Tavš Tet >>>> --Original Message-- > From: xxx@gmail.com > Date: 18.10.2019 21:08 > To: tet@tet.lv > Subject: atbloķēt e-pastu >>> Lūdzu atbloķēt e-pastu xxx@apollo.lv > Visas savas iekārtas esmu iztīrījis no vīrusiem

#### *Attēls 4.3. Datu attīrīšanas piemērs – oriģināla klienta sarakste*

Attēlā (skat. Attēls 4.3.) ir redzami oriģināli, neapstrādāti e-pastu dati. Šajā gadījumā šī ir sarakste starp klientu un uzņēmuma apkalpošanas un palīdzības nodaļas darbinieku.

Lūdzu atbloķēt e-pastu xxx@apollo.lv > Visas savas iekārtas esmu iztīrījis no vīrusiem

#### *Attēls 4.4. Datu attīrīšanas piemērs – e-pastu dati pēc pirmā soļa attīrīšanas*

Attēlā (skat. Attēls 4.4.) ir redzami attīrīti e-pasta dati pēc pirmā attīrīšanas soļa darbības principa.

atbloket pastu visas savas iekartas esmu iztirijis virusiem

#### *Attēls 4.5. Datu attīrīšanas piemērs – e-pastu dati pēc otrā soļa attīrīšanas*

Attēlā (skat. Attēls 4.5.) ir redzami attīrīt e-pasta dati pēc otrā attīrīšanas soļa darbības principiem.

#### 4.3.3 Datu attīrīšanas soļu reāls vizualizējums ar piemēru **negatīviem** datiem.

*Jūtīga informācija tiek aizstāta ar simbolu 'x'.*

Paldies par informāciju, xxx On Mon, Oct 14, 2019 at 8:08 AM Tet <tet@tet.lv> wrote: > Sveiki! >>> Norādītajā adresē varam piedāvāt tikai 4G internetu. Savukārt internets, > ko nodrošinām ar kabeļiem, nav pieejams. >> Sīkāk par 4G internetu var lasīt šeit: > <https://www.tet.lv/internets/pieslegumu-veidi/4g-internets> >>> Ja piedāvājums ir saistošs, lūdzam atsūtīt kontakttālruna numuru, lai > varam vienoties par pakalpojuma pieslēgšanu. >>> Patīkamu dienu vēlot, > Tavs Tet >>> --Original Message-- > From: xxx@gmail.com > Date: 13.10.2019 16:59 > To: tet@tet.lv > Subject: Internets Ventspili >>> Sveiki ! >> Vēlējos noskaidrot par pakalpojumu pieejamību Ventspili. > Kāda veida interneta savienojumus jūs piedāvājat Ventspilī, Ozolu iela x? >> Ar cieņu, > xxx|>

#### **Attēls 4.6. Datu attīrīšanas piemērs – oriģināla klienta sarakste**

Attēlā (skat. Attēls 4.6.) ir redzami oriģināli, neapstrādāti e-pastu dati. Līdzīgi kā iepriekš šī ir sarakste starp klientu un uzņēmuma apkalpošanas un palīdzības nodaļas darbinieku.

Sveiki ! >> Vēlējos noskaidrot par pakalpojumu pieejamību Ventspili. > Kāda veida interneta savienojumus jūs piedāvājat Ventspilī, Ozolu iela x ? >>

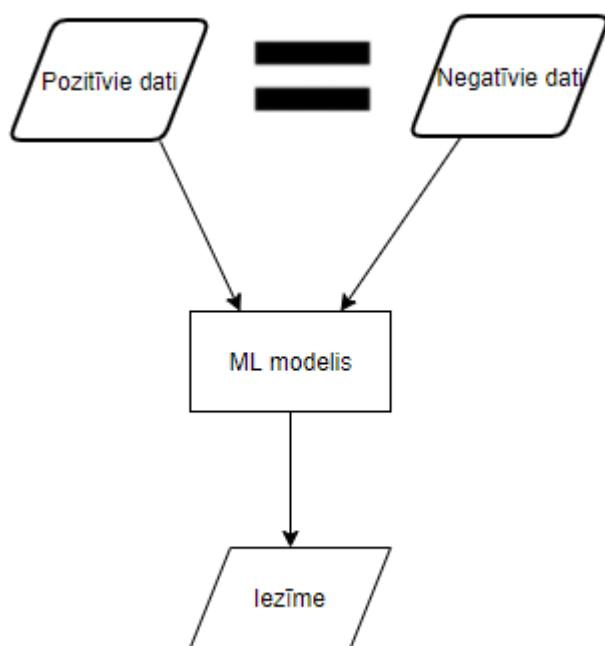
#### **Attēls 4.7. Datu attīrīšanas piemērs – e-pastu dati pēc pirmā soļa attīrīšanas**

Attēlā (skat. Attēls 4.7.) ir redzami attīrīti e-pasta dati pēc pirmā attīrīšanas soļa darbības principa.

velejos noskaidrot pakalpojumu piejamibu ventspili kada veida interneta savienojumus piedavajat ventspili ozolu iela

#### **Attēls 4.8. Datu attīrīšanas piemērs – e-pastu dati pēc otrā soļa attīrīšanas**

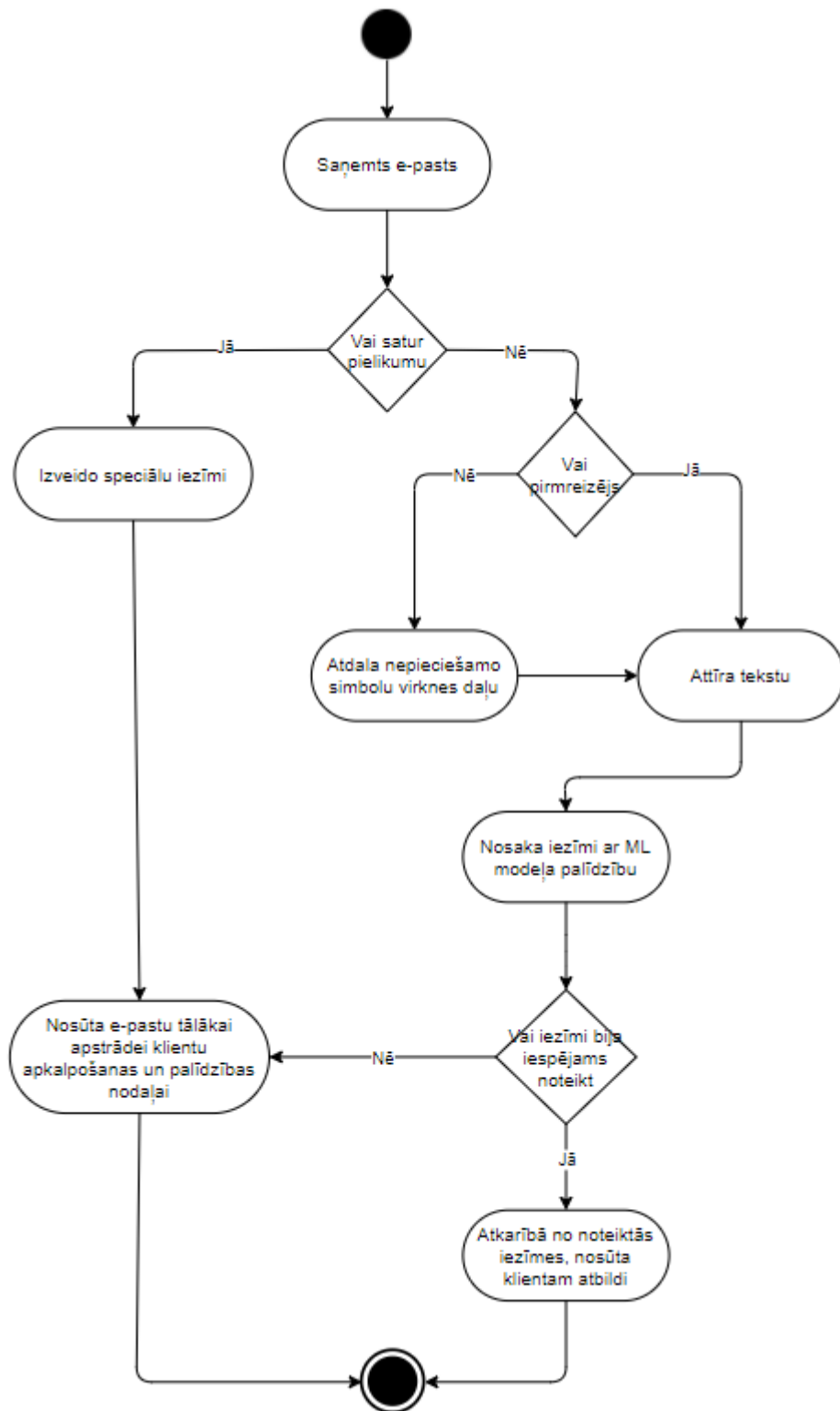
Attēlā (skat. Attēls 4.8.) ir redzami attīrīt e-pasta dati pēc otrā attīrīšanas soļa darbības principiem.



*Attēls 4.9. Triviāls ML modeļa princips*

Attēlā (skat. Attēls 4.9.) Tiek attēlots triviāls princips kā tiek trenēts ML modelis. Modelim tiek padots vienāds skaits pozitīvo datu un negatīvo datu, pēc kā tiek veidota iezīme, kas parāda vai dati pieder pie vienas kopas vai otras.

#### 4.4 Daļējs funkciju projektējums



Attēls 4.10. AEASS aktivitāšu diagramma

Diagramma (skat Attēls 4.10) ir parādītas sistēmas galvenās aktivitātes.

## 5 TESTĒŠANA

AEAASS izstrādes laikā tika veikti testi, lai nodrošinātu sistēmas korektu darbību. Dokumentā tiek aprakstīta sistēmas teksta un e-pastu apstrādes funkcijas.

Tabula 5.1

ID	Ievaddati	Sagaidāmais rezultāts	Reālais rezultāts
isValidEmail	E-pasts tikai ar tekstu	Funkcija atgriež "True"	Izpildās
<b>Apraksts</b>	E-pasts tikai ar pielikumu	Funkcija atgriež "False"	Izpildās
Tiek pārbaudīts e-pasta saturs	E-pasts gan ar pielikumu, gan tekstu	Funkcija atgriež "False"	Izpildās

Tabula 5.2

ID	Ievaddati	Sagaidāmais rezultāts	Reālais rezultāts
IsFirst	E-pasts ar pielikumu	Funkcija atgriež kļūdu	Izpildās
<b>Apraksts</b>	Pirmreizējs e-pasts	Funkcija atgriež "True"	Izpildās
Tiek pārbaudīts vai e-pasts ir pirmreizējs	E-pasts, kas ir daļa no sarakstes ar klientu palīdzības nodaļas operatoru	Funkcija atgriež "False"	Izpildās

Tabula 5.3

ID	Ievaddati	Sagaidāmais rezultāts	Reālais rezultāts
sliceText	Teksts ar sarakstes elementiem	Funkcija atgriež sarakstu ar atrastajiem elementiem, kas liecina par saraksti	Izpildās
<b>Apraksts</b>	Pirmreizējs teksts	Funkcija atgriež tukšu sarakstu	Izpildās
Tiek pārbaudīts vai teksts tiek korekti attīrīts, ja tā ir sarakste			

Tabula 5.4

ID	Ievaddati	Sagaidāmais rezultāts	Reālais rezultāts
cleanText	Simbolu virkne, kas satur tikai dažādus simbolus (\$,*,@, u.tml.)	Funkcija atgriež tukšu simbolu virkni	Izpildās
<b>Apraksts</b>	Simbolu virkne, kas satur burtus ar mīkstinājuma un garumzīmēm	Funkcija atgriež simbolu virkni bez mīkstinājuma vai garumzīmēm	Izpildās
Pārbauda vai teksts tiek korekti attīrīts	Simbolu virkne, kas satur ne-latīnu valodas burtus	Funkcija atgriež tukšu simbolu virkni	Izpildās

Tabula 5.5

ID	Ievaddati	Sagaidāmais rezultāts	Reālais rezultāts
isClassified	Attīrīta simbolu virkne	Funkcija atgriež "True"	Izpildās
<b>Apraksts</b>	Neattīrīta simbolu virkne	Funkcija atgriež "True" (pat ja iezīme nav faktiski korekta)	Izpildās
Pārbauda vai ML modelis korekti piešķir iezīmi	Tukša simbolu virkne	Funkcija atgriež kļūdu	Izpildās

Tabula 5.6

ID	Ievaddati	Sagaidāmais rezultāts	Reālais rezultāts
doRespond	Simbolu virkne ar iezīmi “1”	Funkcija atgriež “1” iezīmes tekstu	Izpildās
<b>Apraksts</b>	Dimbolu virkne ar iezīmi “0”	Funkcija atgriež “2” iezīmes tekstu	Izpildās
Pārbauda, vai pēc iezīmes tiek noteikts korekts atbildes teksts			

## 6 PROJEKTA ORGANIZĀCIJA

Projektu autors veica pilnīgi patstāvīgi, projekta ietvaros autoram bija nepieciešams apgūt sekojošas tehniskas prasmes: programmēšanas valodas python apgūšana, NLP principu apgūšana, ML un neironu tīklu pamatprincipu apgūšana.

Šis ir unikāls un pirmreizējs projekts uzņēmumā, kas tiks turpmāk attīstīts ar autora palīdzību. Lai gan uzņēmums strādā pēc SCRUM metodes, autors darbu veica brīvā stilā, neievērojot konkrētas darba metodes, jo skaitījās kā attālināti piesaistītais darba spēks. Autors pats uzņēmās atbildību veikt komunikāciju ar citām uzņēmuma nodaļām, lai saprastu to vajadzības un intereses projekta ietvaros, autors piedalījās projekta saistītajās tikšanās un konferencēs. Kā galvenais konsultants autoram bija uzņēmuma Lielo Datu apstrādes nodaļas grupas vadītājs.

Uzņēmums patur tiesības par informācijas konfidencialitāti, un autors ir aicināts šīs tiesības ievērot. Uzņēmums atļaujas atklāt savu faktisko nosaukumu, bet autors ir aicināts publicēt un demonstrēt tikai atsevišķas sistēmas pirmkoda daļas.



## **7 KONFIGURĀCIJAS PĀRVALDĪBA**

Sistēma tiek iešūta SMTP servera konfigurācijā, programmkoda rakstīšanai tika izmantota atvērtā pirmkoda programmatūra “Jupyter Notebook”. Tika veikta vietēja versijas kontrole, GIT sistēmā.

## 8 DARBIETILPĪBAS NOVĒRTĒJUMS

Darbietilpība tiek novērtēta balstoties uz “QSM Benchmark Tables Business Systems Implementation Unit (New and Modified IU) Benchmarks” novērtējumu tabulu. Projekts prasīja ļoti daudz plānošanas, par cik autors ko šādu izstrādājis pirmo reizi. Gandrīz viens cilvēkmēnesis tika veltīts projekta plānošanai un komunikācijai. Pirmkoda rindiņu skaits, neskaitot papildus pielāgoto pirmkodu iešūšanai SMTP serverī sastāda 802 koda rindiņas. Jāņem vērā, ka tika izmantota programmēšanas valoda Python un gatavās funkciju bibliotēkas.

Balstoties uz tabulas pirmā kvartīļa informāciju, kas apraksta mazo projektu darbietilpību, ņemot vērā mazāku koda rindiņu skaitu, bet lielāku implementāciju vienību skaitu, tiek secināts, ka projekta izstrādes novērtējums pēc tabulas atbilst 2.5 – 3 personmēnešiem.

## **9 SECINĀJUMI UN AUTORA KOMENTĀRI**

Projekta risinājuma gaitā autors guva unikālas zināšanas gan par uzņēmuma iekšējām darbības struktūrām un principiem, gan par projekta saistītajām tehniskajām lietām. Autors apguva patstāvīgas prasmes projektu izstrādāšanā un problēmu risināšanā. Apgūtās zināšanas ļaus autoram turpināt AEAASS pilnveidošanu un uzturēšanu uzņēmuma interesēs, kā arī ļaus piedalīties līdzīgos izstrādes projektos vai veikt patstāvīgu sistēmas izstrādi.

## 10 IZMANTOTĀ LITERATŪRA UN AVOTI

1. LVS 68:1996 "Programmatūras prasību specifikācijas ceļvedis"
2. "Intrneta agresivitātes indekss" [Tiešsaiste]. Pieejams:  
<http://barometrs.korpuss.lv/?from=2013-10-13&to=2014-10-13&site=lv&section=stopwords>
3. Latvijas Nacionālais standartizācijas un metroloģijas centrs, Latvijas standarts LVS 72:1996, Ieteicamā prakse programmatūras projektējuma aprakstīšanai
4. "Your Guide to Natural Language Processing (NLP)" [Tiešsaiste]. Pieejams:  
<https://towardsdatascience.com/your-guide-to-natural-language-processing-nlp-48ea2511f6e1>
5. "LSTM Layer" [Tiešsaiste]. Pieejams: [https://keras.io/api/layers/recurrent\\_layers/lstm/](https://keras.io/api/layers/recurrent_layers/lstm/)

## PIELIKUMI

### Teksta apstrādes un ML modeļa trenēšanas pirmkods

```
import csv
import sys
Import os
import re
import enchant
import nltk
import nltk.corpus
from nltk.tokenize import word_tokenize
import numpy as np
import pandas as pd
from pandas import ExcelWriter
from pandas import ExcelFile

pd.options.mode.chained_assignment = None
maxInt = sys.maxsize
maxInt = int(maxInt/10000000000)
csv.field_size_limit(maxInt)

d1 = pd.read_csv(r'E:\email_documents\email_files\10-440-400.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d2 = pd.read_csv(r'E:\email_documents\email_files\11-400-360.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d3 = pd.read_csv(r'E:\email_documents\email_files\12-360-320.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d4 = pd.read_csv(r'E:\email_documents\email_files\13-320-280.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d5 = pd.read_csv(r'E:\email_documents\email_files\14-280-240.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d6 = pd.read_csv(r'E:\email_documents\email_files\15-240-200.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d7 = pd.read_csv(r'E:\email_documents\email_files\16-200-150.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d8 = pd.read_csv(r'E:\email_documents\email_files\17-150-100.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d9 = pd.read_csv(r'E:\email_documents\email_files\18-100-50.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d10 = pd.read_csv(r'E:\email_documents\email_files\19-50-0-END.csv',
sep=';', header=None, engine='python', encoding="utf-8-sig",
error_bad_lines=False, usecols=[17,24])
d11 = pd.read_csv(r'E:\email_documents\email_files\3-720-700.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d12 = pd.read_csv(r'E:\email_documents\email_files\4-700-660.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
```

```

d13 = pd.read_csv(r'E:\email_documents\email_files\5-660-620.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d14 = pd.read_csv(r'E:\email_documents\email_files\6-620-580.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d15 = pd.read_csv(r'E:\email_documents\email_files\7-580-530.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d16 = pd.read_csv(r'E:\email_documents\email_files\8-530-480.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])
d17 = pd.read_csv(r'E:\email_documents\email_files\9-480-440.csv', sep=';',
header=None, engine='python', encoding="utf-8-sig", error_bad_lines=False,
usecols=[17,24])

df=pd.concat([d1,d2,d3,d4,d5,d6,d7,d8,d9,d10,d11,d12,d13,d14,d15,d16,d17])
df.dropna(subset = [17,24], inplace=True)

#Drop invalid rows (in email column - drop invalid emails that do not
contain @)
df = df[df[17].str.contains('@')]
df=df[~df[17].str.contains('tetteikals')]

#reset column and row indexes
df.reset_index(inplace=True, drop=True)
df.columns = range(df.shape[1])

#First part of cleaning
#This part finds tet questions from templates phrases in dataframe

#initialize set of template questions or phrases
template_list = ['klienta numuru vai personas kodu', 'klienta numuru vai >
personas kodu',
'nepieciešami klienta dati', 'personas kodu vai klienta numuru', 'norādīt
klienta numuru', 'personas kodu un/vai klienta numuru',
'klienta numuru vai personas kodu', 'klienta līguma numuru', 'precizējiet
personas kodu', 'personas kodu un/vai klienta > numuru',
'abonenta personas kodu', 'precizēt klienta datus, līguma numuru', 'klienta
> līguma numuru', 'klienta līguma numuru', 'klienta > līguma numuru']
df = df[df[1].str.contains('|'.join(template_list), case =True)]

#Leave only most popular private emails
private_list = ['gmail.com', 'inbox.lv', 'mail.ru', 'bk.ru', 'apollo.lv',
'hotmail.com', 'zb.lv', 'yahoo.com']
df = df[df[1].str.contains('|'.join(private_list), case=True)]

#Select emails only from clients
l = ['tet', 'lattelecom']
df = df[~df[0].str.contains('|'.join(l), case = True)]
df = df[~df[1].str.contains('Tet E-veikalā')]
df = df[~df[1].str.contains('Tet veikals')]
df = df[~df[1].str.contains('lattelecomveikals')]
df = df[~df[1].str.contains('Lai mēs varētu sameklēt un novirzīt maksājuma
uzdevumā norādīto naudas summu')]
df = df[~df[1].str.contains('nepieciešams atsūtīt ārsta apstiprinātu
izziņu')]
df = df[~df[1].str.contains('nav pielikuma')]

```



```

final_clean[1] = final_clean[1].str.replace('\w*inbox', '')

#All to lower
final_clean[1] = final_clean[1].str.lower()

#remove all kinds of symbols
final_clean[1] =
final_clean[1].str.replace('ä', 'a').replace('à', 'a').replace('ä',
'a').replace('á', 'a').replace('â', 'a').replace('ã', 'a').replace('å', 'a')
final_clean[1] = final_clean[1].str.replace('ë', 'e').replace('é',
'e').replace('è', 'e').replace('ê', 'e').replace('ë', 'e')
final_clean[1] = final_clean[1].str.replace('ś', 's')
final_clean[1] = final_clean[1].str.replace('š', 's')
final_clean[1] = final_clean[1].str.replace('ï', 'i').replace('ì',
'i').replace('í', 'i').replace('î', 'i').replace('ï', 'i')
final_clean[1] = final_clean[1].str.replace('ü', 'u')
final_clean[1] = final_clean[1].str.replace('č', 'c')
final_clean[1] = final_clean[1].str.replace('ñ', 'n').replace('ñ', 'n')
final_clean[1] = final_clean[1].str.replace('ķ', 'k')
final_clean[1] = final_clean[1].str.replace('ġ', 'g')
final_clean[1] = final_clean[1].str.replace('ĺ', 'l')
final_clean[1] = final_clean[1].str.replace('ļ', 'l')
final_clean[1] = final_clean[1].str.replace('ž', 'z')
final_clean[1] = final_clean[1].str.replace('ù', 'u').replace('ú',
'u').replace('ù', 'u').replace('û', 'u').replace('ü', 'u')
final_clean[1] = final_clean[1].str.replace('☺', '')
final_clean[1] = final_clean[1].str.replace('☹', '')
final_clean[1] = final_clean[1].str.replace('ee', 'e')
final_clean[1] = final_clean[1].str.replace('aa', 'a')
final_clean[1] = final_clean[1].str.replace('ii', 'i')
final_clean[1] = final_clean[1].str.replace('uu', 'u')
final_clean[1] = final_clean[1].str.replace('ee', 'e')
final_clean[1] = final_clean[1].str.replace('sh', 's')
final_clean[1] = final_clean[1].str.replace('lj', 'l')
final_clean[1] = final_clean[1].str.replace(' cee', ' ')
final_clean[1] = final_clean[1].str.replace(' mb ', ' ')
final_clean[1] = final_clean[1].str.replace(' uz ', ' ')
final_clean[1] = final_clean[1].str.replace(' no ', ' ')
final_clean[1] = final_clean[1].str.replace(' ne ', ' ')
final_clean[1] = final_clean[1].str.replace(' ja ', ' ')
final_clean[1] = final_clean[1].str.replace(' ko ', ' ')
final_clean[1] = final_clean[1].str.replace(' ap ', ' ')
final_clean[1] = final_clean[1].str.replace(' ir ', ' ')
final_clean[1] = final_clean[1].str.replace(' un ', ' ')
final_clean[1] = final_clean[1].str.replace(' ka ', ' ')
final_clean[1] = final_clean[1].str.replace(' bet ', ' ')
final_clean[1] = final_clean[1].str.replace(' ar ', ' ')
final_clean[1] = final_clean[1].str.replace(' par ', ' ')
final_clean[1] = final_clean[1].str.replace(' es ', ' ')
final_clean[1] = final_clean[1].str.replace(' nu ', ' ')
final_clean[1] = final_clean[1].str.replace(' tas ', ' ')
final_clean[1] = final_clean[1].str.replace(' man ', ' ')
final_clean[1] = final_clean[1].str.replace(' to ', ' ')
final_clean[1] = final_clean[1].str.replace(' kas ', ' ')
final_clean[1] = final_clean[1].str.replace(' tad ', ' ')
final_clean[1] = final_clean[1].str.replace(' ko ', ' ')
final_clean[1] = final_clean[1].str.replace(' vai ', ' ')
final_clean[1] = final_clean[1].str.replace(' ta ', ' ')
final_clean[1] = final_clean[1].str.replace(' ari ', ' ')
final_clean[1] = final_clean[1].str.replace(' tu ', ' ')
final_clean[1] = final_clean[1].str.replace(' gan', ' ')

```



[illegible]

```

final_clean[1] = final_clean[1].str.replace(':', ' ')
final_clean[1] = final_clean[1].str.replace('"', ' ')
final_clean[1] = final_clean[1].str.replace('€', ' ')
final_clean[1] = final_clean[1].str.replace(' lv', ' ')
final_clean[1] = final_clean[1].str.replace('lv ', ' ')
final_clean[1] = final_clean[1].str.replace(' lv ', ' ')
final_clean[1] = final_clean[1].str.replace('labdien', ' ')
final_clean[1] = final_clean[1].str.replace('labvakar', ' ')
final_clean[1] = final_clean[1].str.replace('sveiki', ' ')
final_clean[1] = final_clean[1].str.replace('ludzu', ' ')
final_clean[1] = final_clean[1].str.replace('paldies', ' ')

#remove more stuff
final_clean[1] = final_clean[1].str.replace('iphone', ' ')
final_clean[1] = final_clean[1].str.replace('twitpd', ' ')
final_clean[1] = final_clean[1].str.replace('from', ' ')
final_clean[1] = final_clean[1].str.replace('reply message', ' ')
final_clean[1] = final_clean[1].str.replace('lattelecom', ' ')
final_clean[1] = final_clean[1].str.replace('tet', ' ')
final_clean[1] = final_clean[1].str.replace('sent from my iphone', ' ')
final_clean[1] = final_clean[1].str.replace('no subject', ' ')
final_clean[1] = final_clean[1].str.replace('from tet sent friday may pm
to', ' ')
final_clean[1] = final_clean[1].str.replace('tavs bezmaksas pasts inbox', ' ')
final_clean[1] = final_clean[1].str.replace('uga records media
wwwugarecordscom', ' ')

#remove non-latin letters
final_clean[1] = final_clean[1].str.replace('[^\x00-\x7F]+', '')

#remove numbers and single chars and whitespaces
final_clean[1] = final_clean[1].str.replace('[0-9]', '')
final_clean[1] = final_clean[1].str.replace(' +', ' ')
#exception case
final_clean[1] = final_clean[1].str.replace('sent mail.*', ' ')
final_clean[1] = final_clean[1].str.replace('sent my.*', ' ')
final_clean[1] = final_clean[1].str.replace('image.*', ' ')
final_clean[1] = final_clean[1].str.replace('ar cienu.+', ' ')
final_clean[1] = final_clean[1].str.replace('reply to mess.+', ' ')
final_clean[1] = final_clean[1].str.replace('sent from.+', ' ')
final_clean[1] = final_clean[1].str.replace('am deply.*', ' ')
final_clean[1] = final_clean[1].str.replace('sincerely.+', ' ')
final_clean[1] = final_clean[1].str.replace('forwarded.+', ' ')
final_clean[1] = final_clean[1].str.replace(' lpojums', ' pakalpojums')
final_clean[1] = final_clean[1].str.replace(' elotu', ' novelotu')
final_clean[1] = final_clean[1].str.replace(' lpojuma', ' pakalpojuma')
final_clean[1] = final_clean[1].str.replace(' lpojumas', ' pakalpojumas')
final_clean[1] = final_clean[1].str.replace(' lpojumi', ' pakalpojumi')
final_clean[1] = final_clean[1].str.replace(' lpojumiem', ' pakalpojumiem')
final_clean[1] = final_clean[1].str.replace(' lpojums', ' pakalpojums')
final_clean[1] = final_clean[1].str.replace(' lpojumu', ' pakalpojumu')
final_clean[1] = final_clean[1].str.replace(' lpojums', ' pakalpojums')

#remove template text
final_clean[1] = final_clean[1].str.replace('^jautajums.+jautajums', ' ')

#remove more numbers and whitespaces
final_clean[1] = final_clean[1].str.replace('^\s+|\s+$', '')

```

```

final_clean[1] = final_clean[1].str.replace('\s\s+', ' ')

#remove www
final_clean[1] = final_clean[1].str.replace('\bwww.*', ' ')
final_clean[1] = final_clean[1].str.replace('\swww.*', ' ')

#remove lone characters
final_clean[1] = final_clean[1].str.replace('(^| ).( |$)', ' ')

#remove empty rows (again after cleaning)
final_clean[1].replace('', np.nan, inplace=True) #replace empty strings
withn NaN
final_clean.dropna(subset=[1], inplace=True)
#remove rows where word count is small
count = final_clean[1].str.split().str.len()
~(count == 1)
final_clean = final_clean[~(count == 1)]
final_clean = final_clean[~(count == 2)]
final_clean = final_clean[~(count == 3)]
final_clean = final_clean[~(count == 4)]
final_clean = final_clean[~(count == 5)]

#Set max word count for data
max_size = 25
final_clean[1] =
final_clean[1].str.split(n=max_size).str[:max_size].str.join(' ')

# # ML PART

##### typ2 #####
from sklearn.feature_extraction.text import CountVectorizer
from sklearn import metrics
vect = CountVectorizer()
from tensorflow.keras.preprocessing import sequence
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Embedding
from tensorflow.keras.layers import LSTM
from tensorflow.keras.utils import to_categorical
from tensorflow.keras.preprocessing.text import Tokenizer
from sklearn.metrics import accuracy_score
from tensorflow.keras.preprocessing.sequence import pad_sequences
#####
import warnings
warnings.filterwarnings('ignore')

#Load labeled dataframe
df = pd.read_excel(r'E:\email_documents\email\labeled_data.xlsx')
test_df = pd.read_excel(r'E:\email_documents\email\test.xlsx')
df['text'] = df['text'].astype('str')
test_df['text'] = test_df['text'].astype('str')

#Model data
X_train = df['text']
X_test = test_df['text']
y_train = df.label
y_test = test_df.label
tokenize = Tokenizer()

```

```

tokenize.fit_on_texts(X_train.values)

#tokenize data
X_train = tokenize.texts_to_sequences(X_train)
X_test = tokenize.texts_to_sequences(X_test)

#make tokenized data into sequence
max_lenght = max([len(s.split()) for s in df['text']]) #length of input
sequences to Embedding layer
X_train = pad_sequences(X_train, max_lenght)
X_test = pad_sequences(X_test, max_lenght)

#size of the output vectors (bigger num => more params)
EMBEDDING_DIM = 128

#vocabulary size
unknown = len(tokenize.word_index)+1

model = Sequential()
model.add(Embedding(unknown, EMBEDDING_DIM, input_length=max_lenght))
model.add(LSTM(units=10, dropout=0.1, recurrent_dropout=0.1 ))
model.add(Dense(1, activation='sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='adam',
metrics=['accuracy'])
print(model.summary())
model.fit(np.asarray(X_train), np.asarray(y_train), batch_size=64,
epochs=8, verbose=1)

final_pred = model.predict_classes(X_test) #predict_proba
print(accuracy_score(y_test, final_pred))
y_pred_df = pd.DataFrame(final_pred)
y_pred_df.columns = ['prediction']
final_test = pd.concat([test_df, y_pred_df], axis=1)
final_test.to_excel(r'E:\email_documents\email\result_deep_learning.xlsx')

max_length = 25

a=b = ['man ir neapmaksats rekins nakamnedel apmaksasu']

a = tokenize.texts_to_sequences(a)
a = pad_sequences(a, max_length)
print(model.predict_classes(a))

```

Kvalifikācijas darbs „*Automātiska e-pastu apstrādes un atbildes sniegšanas sistēma*”  
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: **Edvards Bukovskis** \_\_\_\_\_ **09.01.2021.**

Rekomendēju darbu aizstāvēšanai

Darba vadītājs/a: **Prof. Jānis Zuters** \_\_\_\_\_ **.01.2021.**

Recenzents: **Baltcom SIA, struktūrvienības vadītājs, Lauris Raipulis**

Darbs iesniegts 11.01.2021.

Kvalifikācijas darbu pārbaudījumu komisijas sekretāre: **Darja Solodovņikova** \_\_\_\_\_

Darbs aizstāvēts kvalifikācijas darbu pārbaudījuma komisijas sēdē

\_\_\_\_.01.2021. prot. Nr. \_\_\_\_\_

Komisijas sekretārs(-e): \_\_\_\_\_