

VŠB – Technická univerzita Ostrava  
Fakulta elektrotechniky a informatiky  
Katedra informatiky

# **Analýza emailové komunikace**

# **Analysis of Email Communication**

2018

Veronika Uhrová

## Zadání diplomové práce

Student:

**Bc. Veronika Uhrová**

Studijní program:

N2647 Informační a komunikační technologie

Studijní obor:

2612T025 Informatika a výpočetní technika

Téma:

Analýza emailové komunikace

Analysis of Email Communication

Jazyk vypracování:

čeština

Zásady pro vypracování:

Cílem práce je návrh a implementace systému na analýzu emailové komunikace a vizualizaci výstupů. Systém bude pracovat s reálnými daty a budou navrženy, popsány a vyhodnoceny experimenty s těmito daty. Pro implementaci je doporučen jazyk C#.

1. Rešerše obdobných řešení a analytických přístupů.
2. Návrh a implementace metody na získávání emailových zpráv z vybraného zdroje.
3. Výběr a implementace metod strojového učení a analýzy sítí vhodných pro analýzu emailové komunikace.
4. Návrh a implementace uživatelského rozhraní na analýzu emailové komunikace a vizualizaci analytických výstupů.
5. Dokumentovaná implementace systému.

Seznam doporučené odborné literatury:

- [1] S. Zehnalova, Z. Horak, M. Kudělka. Email Conversation Network Analysis: Work Groups and Teams in Organizations. ASONAM 2015.  
[2] Tang, G., Pei, J., Luk, W. S. (2014). Email mining: tasks, common techniques, and tools. Knowledge and Information Systems, 41(1), 1-31.

Další podle pokynů vedoucího práce.

Formální náležitosti a rozsah diplomové práce stanoví pokyny pro vypracování zveřejněné na webových stránkách fakulty.

Vedoucí diplomové práce: **doc. Mgr. Miloš Kudělka, Ph.D.**

Datum zadání: 01.09.2017

Datum odevzdání: 30.04.2018

---

doc. Ing. Jan Platoš, Ph.D.  
*vedoucí katedry*



---

prof. Ing. Pavel Brandstetter, CSc.  
*děkan fakulty*

Prehlasujem, že som túto prácu vypracovala samostatne. Uviedla som všetky literárne pramene a publikácie, z ktorých som čerpala.

V Ostrave, 27.4.2018

.....

Moje podakovanie patrí predovšetkým doc. Milošovi Kudělkovi, Ph.D. za odborné konzultácie a vedenie mojej diplomovej práce.

## **Abstrakt**

Práca študuje aktuálne metódy pre analýzu emailov a detekciu sociálnych rolí v emailových dátach. Nasleduje zoznámenie sa s emailom a jeho popularitou v súčasnosti. Taktiež práca uvádzajú základné teoretické pojmy a teoretický náhľad na reprezentáciu siete. Uvádzajú sa tu aj základy analýzy sociálnych sietí a detekcie komunit. Ďalej sa tu píše o framework pre detekciu štrukturálnych rolí a ich identifikácie. Na základe týchto poznatkov je vytvorená aplikácia pre analýzu a vizualizáciu analytických výstupov. Na záver sú uvedené prevedené experimenty týkajúce sa identifikácie sociálnych rolí v emailových dátach.

**Klúčové slová:** email, sociálna sieť, sociálna rola, vizualizácia, brokerage

## **Abstract**

This paper studies current methods for analysing emails and detecting social roles in email data. This is followed by getting acquainted with the email and its popularity nowadays. Also, this thesis presents basic theoretical concepts and a theoretical overview of network representation. Here are also the basics of social networking and community detection. There is also written a framework for structural social roles detection and their identification. Based on this knowledge, an application is developed to analyze and visualize analytical outputs. Finally, experiments on the findings of the emailed data are presented.

**Key Words:** email, social network, social role, visualization, brokerage

# Obsah

<b>Zoznam použitých skratiek a symbolov</b>	<b>10</b>
<b>Zoznam obrázkov</b>	<b>11</b>
<b>Zoznam tabuliek</b>	<b>13</b>
<b>1 Úvod</b>	<b>14</b>
1.1 Motivácia . . . . .	14
1.2 Vízia . . . . .	14
1.3 Štruktúra práce . . . . .	15
<b>2 Súvisiace práce</b>	<b>16</b>
<b>3 Emailová komunikácia</b>	<b>18</b>
3.1 Stručná história emailu . . . . .	18
3.2 Štruktúra emailu . . . . .	18
3.3 Emaily v súčasnosti . . . . .	19
<b>4 Definície a klasifikácie</b>	<b>20</b>
4.1 Graf . . . . .	20
4.2 Súvislost grafu . . . . .	21
4.3 Úplný graf . . . . .	21
4.4 Stupeň vrcholu . . . . .	21
4.5 Cesta . . . . .	22
4.6 Uzavretá cesta . . . . .	22
4.7 Komponenta grafu . . . . .	22
4.8 Metriky . . . . .	22
4.8.1 Closeness centrality (Centralita blízkosti) . . . . .	22
4.8.2 Betweeness centrality (Centralita medziľahosti) . . . . .	22
4.8.3 Modularita . . . . .	23
<b>5 Sociálna sieť</b>	<b>24</b>
5.1 História sociálnych sietí . . . . .	24
5.2 Analýza sociálnych sietí . . . . .	24
5.3 Komunity v sociálnych sietach . . . . .	25
5.4 Detekcia komunít . . . . .	26
5.4.1 Louvainov algoritmus pre detekciu komunít . . . . .	26
5.5 Ego siet . . . . .	27
5.5.1 Konštrukcia ego siete . . . . .	27

<b>6 Metódy analýzy sociálnych sietí</b>	<b>29</b>
6.1 SSRM - Framework pre detekciu štrukturálnych rolí v sociálnych sietach . . . . .	29
6.1.1 Rola v kontexte SSRM . . . . .	29
6.1.2 Roly definované v SSRM . . . . .	29
6.1.2.1 Leader . . . . .	30
6.1.2.2 Outermost . . . . .	30
6.1.2.3 Mediator . . . . .	30
6.1.2.4 Outsider . . . . .	30
6.2 Identifikácia štrukturálnych sociálnych rolí . . . . .	30
6.2.0.1 Outsider . . . . .	31
6.2.1 Leader . . . . .	31
6.2.1.1 Closeness centrality (Centralita blízkosti) . . . . .	31
6.2.2 Outermost . . . . .	31
6.2.3 Mediator . . . . .	31
6.2.3.1 LBeweeness . . . . .	31
6.2.3.2 CBetweenness . . . . .	32
6.2.3.3 Normalizovaná verzia CBetweenness . . . . .	32
6.2.3.4 Skóre rozmanitosti . . . . .	33
6.3 Brokerage roly . . . . .	33
6.3.1 Liaison . . . . .	35
6.3.2 Itinerant . . . . .	35
6.3.3 Coordinator . . . . .	35
6.3.4 Gatekeeper . . . . .	36
6.3.5 Representative . . . . .	36
6.3.6 Identifikácia brokerage rolí . . . . .	36
6.4 Analýza ega . . . . .	37
6.4.1 Veľkosť ego siete . . . . .	37
6.4.2 Kompozícia ego siete . . . . .	38
6.4.3 Štruktúra ego siete . . . . .	38
6.4.3.1 Efektívna veľkosť . . . . .	39
<b>7 Aplikácia</b>	<b>40</b>
7.1 Špecifikácia . . . . .	40
7.1.1 Funkčné požiadavky . . . . .	40
7.2 Návrh . . . . .	41
7.2.1 Návrhové vzory . . . . .	42
7.3 Dôležité rozhodnutia . . . . .	44
7.3.1 Dostupnosť dát . . . . .	44
7.3.2 Webová vs. desktopová aplikácia . . . . .	44

7.4	Použité knižnice . . . . .	44
7.5	Import dát . . . . .	46
7.6	Implementácia . . . . .	46
7.6.1	Metóda pre získanie emailových dát . . . . .	46
7.6.2	Konštrukcia siete . . . . .	47
7.6.3	Triedy pre graf, vrcholy a hrany . . . . .	47
<b>8</b>	<b>Experimenty</b>	<b>48</b>
8.1	Analýza tímu . . . . .	48
8.1.1	Príprava a import dát . . . . .	48
8.1.2	Vizualizácia datasetu . . . . .	49
8.1.3	Detekcia komunít . . . . .	50
8.1.3.1	Zmeny komunít v čase . . . . .	52
8.1.4	Ego siet <sup>†</sup> . . . . .	53
8.1.5	Analýza rolí . . . . .	53
8.1.5.1	SSRM . . . . .	53
8.1.5.2	Brokerage . . . . .	54
8.2	Analýza jednotlivca . . . . .	55
8.2.1	Príprava a import dát . . . . .	55
8.2.2	Informácie o datasete . . . . .	56
8.2.3	Detekcia komunít . . . . .	57
8.2.3.1	Zmeny komunít v čase . . . . .	58
8.2.4	Ego siet <sup>†</sup> . . . . .	59
8.2.5	Analýza rolí . . . . .	59
8.2.5.1	SSRM . . . . .	60
8.2.5.2	Brokerage . . . . .	60
<b>9</b>	<b>Záver</b>	<b>62</b>
9.1	Možnosti rozšírenia a zdokonalenia práce . . . . .	62
<b>Literatura</b>		<b>63</b>

## **Zoznam použitých skratiek a symbolov**

MUA	– Mail User Agent
MTA	– Mail Transfer Agent
IMAP	– Internet Message Access Protocol
XML	– eXtensible Markup Language
SSRM	– Structural social role mining framework
SNA	– Social network analysis

## Zoznam obrázkov

1	Akú formu komunikácie preferujete na formálnu komunikáciu? . . . . .	19
2	Neorientovaný graf . . . . .	20
3	Orientovaný graf . . . . .	20
4	Súvislý (1) a nesúvislý graf (2) . . . . .	21
5	Úplný graf . . . . .	21
6	Graf v tvare hviezdy . . . . .	23
7	Sieť s viacerými komunitami . . . . .	25
8	Vizualizácia krokov Louvainovho algoritmu. . . . .	27
9	Príklad ego siete. . . . .	28
10	Príklad brokerage procesu . . . . .	34
11	Liaison brokerage . . . . .	35
12	Itinerant brokerage . . . . .	35
13	Coordinator brokerage . . . . .	35
14	Gatekeeper brokerage . . . . .	36
15	Representative brokerage . . . . .	36
16	Velkosť ega - stupeň uzla: 6 . . . . .	37
17	Málo štrukturálnych dier vs. veľa štrukturálnych dier. . . . .	38
18	Príklad výpočtu efektívnej sily . . . . .	39
19	UseCase Diagram . . . . .	41
20	Diagram komponent znázorňujúci jednotlivé komponenty architektúry aplikácie .	42
21	Triedny diagram - Repository pattern . . . . .	43
22	Model-View-Controller . . . . .	44
23	Jednoduchá sieť vytvorená s použitím knižnice vis.js . . . . .	45
24	Príklad použitia knižnice vis.js . . . . .	45
25	Doménový model . . . . .	46
26	Príklad konfigurácie emailu pre získanie emailov . . . . .	47
27	Základné informácie o tímovej sieti. . . . .	49
28	Najviac používané emailové domény. . . . .	49
29	Vizualizácia siete. . . . .	50
30	Vizualizácia komunít v tímovej sieti za celkový čas . . . . .	51
31	Rozloženie komunít v tímovej sieti za celkový čas . . . . .	51
32	Rozloženie komunít za prvý časový úsek . . . . .	52
33	Rozloženie komunít za druhý časový úsek . . . . .	52
34	Rozloženie komunít za tretí časový úsek . . . . .	53
35	Počet detekovaných štrukturálnych rôl . . . . .	54
36	Desať aktérov s najväčším <i>brokerage</i> skóre . . . . .	54
37	Desať aktérov s najväčším <i>brokerage</i> skóre - graf . . . . .	55

38	Analýza jednotlivca - základná vizualizácia . . . . .	56
39	Analýza jednotlivca - základné štatistiky . . . . .	57
40	Analýza jednotlivca - vizualizácia komunít . . . . .	58
41	Analýza jednotlivca - Vizualizácia komunít v prvom časovom intervale . . . . .	58
42	Analýza jednotlivca - vizualizácia komunít v druhom časovom intervale . . . . .	59
43	Analýza jednotlivca - detail detekcie <i>leader</i> roly . . . . .	60
44	Desať aktérov s najväčším <i>brokerage</i> skórom . . . . .	60
45	Desať aktérov s najväčším <i>brokerage</i> skórom - graf . . . . .	61

## **Zoznam tabuliek**

1	Základné informácie o datasete . . . . .	48
2	Informácie o členoch tímu . . . . .	50
3	Informácie o vytvorennej ego sieti . . . . .	53
4	Informácie o vytvorennej ego sieti . . . . .	59

# 1 Úvod

V stručnom úvode je popísaná motivácia, ktorá viedla k vypracovaniu tejto diplomovej práce a vízia toho, čo sa malo dosiahnuť a hrubá štruktúra vypracovaného textu.

## 1.1 Motivácia

S cieľom uľahčiť používanie emailov a prebádať podnikateľský potenciál emailov, analýza emailov dosiahla pozoruhodný pokrok nielen v oblasti výskumu, ale aj v praxi. Emaily možno považovať za zmiešanú štruktúru obsahujúcu údaje o ľuďoch zo sociálnych alebo aj organizačných aspektov.

### **Obsah emailu ako textové a netextové dát**

Emaily sú písané viac stručne ako väčšina ostatných dokumentov, často obsahujú hovorové výrazy a abreviácie, ktoré sa nenachádzajú v bežných slovníkoch, preto štandardné techniky analýzy textov pri práci s emailovými dátami nemusia byť efektívne.

Emaily tiež obsahujú bohatšie typy dát, ako napríklad URL linky, HTML tagy alebo obrázky. Niektoré štúdie jednoducho zjednodušia tieto netextové dátové vstupy v štádiu predpripravovania dát - vymažu ich a ďalej pracujú len s textovými dátami. Tieto netextové dátá však môžu byť užitočné v iných oblastiach, ako napríklad detekcia spamu.

### **Emaily reprezentujúce ľudské sociálne organizačné vzťahy**

Emailová aktivita sama o sebe reprezentuje bohaté ľudské sociálne a organizačné vzťahy, ktoré spájajú ľudí do komunít a komplexných systémov. Porozumenie organizačných štruktúr alebo vzťahov naprieč ľuďmi v organizácii môže byť veľmi užitočné aj v reálnom živote. Hlavné problémy, ktoré sú investigované v analýze emailov sú detekcia spamu, kategorizácia emailov, analýza kontaktov, analýza vlastností emailových sietí a vizualizácia emailov.

## 1.2 Vízia

Cieľom práce je oboznámiť čitateľa s oblasťou sociálnych sietí a špeciálne s tému analýzy emailových dát a tieto znalosti demonštrovať nad reálnymi emailovými dátami. Pre uskutočnenie tohto cieľa je potrebné naštudovať informácie z oblasti analýzy emailov, reprezentácie emailu v sieti a vizualizácie sociálnych sietí vrátane aktuálnych metód publikovaných v článkoch. K tomu sa viaže tiež prieskum reprezentácie a konštrukcie emailu ako prvku sociálnej siete.

Ďalej boli vybrané metódy detekcie rolí v sociálnej sieti a navrhnutá aplikácia, ktorá umožňuje analyzovať a vizualizovať analytické výsledky. V tejto aplikácii s jednoduchým a použiteľným užívateľským rozhraním sú implementované vybrané metódy analýzy a je navrhnutá prehľadná vizualizácia vzťahov. Nakoniec je vytvorená analýza tímu podľa emailových dát a porovnanie dvoch prvkoch siete a výsledky experimentov sú zrozumiteľne prezentované.

### **1.3 Štruktúra práce**

V prvej kapitole je uvedený prieskum o aktuálnych vedeckých článkoch, ktoré sa zaoberejú analýzou emailov a reprezentáciou emailu v sociálnych sieťach. Ďalej sa čitateľ zoznámi s emailom ako komunikačným prostriedkom a dozvie sa, ako sú na tom emaily s popularitou aktuálne. Potom je uvedený stručný prehľad teórie grafov a definícií určitých pojmov, ktorý je nevyhnutný k porozumeniu ďalších kapitol. V ďalšej kapitole píšem o sociálnych sieťach, ich histórii a analýze sociálnych sietí, komunitnej štruktúre sociálnych sietí a ego sietach. Neskôr prechádzam k popisu a reprezentácii frameworku pre detekciu štrukturálnych rolí, popisujem sociálne roly definované v rámci tohto frameworku a následne v ďalšej kapitole referujem pomocou akých metód sa sociálne roly v rámci tohto frameworku identifikujú. Ďalej popisujem ďalšiu metódu pre identifikáciu rolí zo sociálnych sietí - *brokerage*. Na základe všetkých poznatkov práce je navrhnutá aplikácia vhodná k sledovaniu výsledkov navrhnutých metód pre analýzu emailových sietí. Ešte pred záverom sú uvedené výsledky prevedených experimentov týkajúcich sa poznatkami skúmanej sociálnej siete.

## 2 Súvisiace práce

Pre odhalovanie vzťahov medzi ľuďmi, skupinami a organizáciami z emalových sietí boli aplikované mnohé techniky a modely analýzy sociálnych sietí. Mnoho štúdií použilo maily spoločnosti *Enron* kvôli nedostatku dostupných veľkých súborov.

Napríklad Diesner, Carley a Frantz v [1] zkonštruovali z mailovej komunikácie spoločnosti Enron orientovaný graf zo vzťahu odosielateľ-príjemca, kde hrany boli vážené frekvenciou mailov, ktoré si medzi sebou poslali v čase. Potom aplikovali techniky analýzy sociálnych sietí. V práci popísali, ako vylepšili originálnu sadu a súčasné zistenia ich investigáciou vďaka analýze sociálnych sietí. Skúmajú dynamiku, štruktúru a vlastnosti organizačnej komunikačnej siete ako aj charakteristiky a vzory komunikačného správania zamestnancov z rôznych organizačných levelov. Zistili, že počas obdobia krízy sa komunikácia medzi zamestnancami stala viac rôznorodejšia v súvislosti so zavedenými kontaktami a formálnymi rolami. Taktiež počas obdobia kríz, predtým nekomunikujúci zamestnanci sa začali zapájať do vzájomného rozhovoru, takže interpersonálna komunikácia bola intenzívnejšia a sieť sa tým rozširovala. Tieto zistenia poskytli cenný pohľad do organizačnej krízy reálneho sveta, čo môže byť ďalej využité pre validáciu alebo tvorbu teórií a dynamických modelov organizačných kríz a tým to vedie k lepšiemu porozumeniu základných príčin organizačných kríz v organizáciách.

Xiaoyan Fu v [2] prezentoval rôzne metódy pre vizualizáciu emailových sietí. Vizualizácia objavuje komunikačné vzory medzi rôznymi skupinami, zobrazuje centrálnu analýzu s dôrazom na významné uzly. V práci zkonštruovali 2D vizualizáciu temporálnej emailovej siete, ktorá analyzuje vývoj emailových vzťahov, ktoré sa menia v priebehu času a zobrazenie prostredia pre nájdenie sociálnych kruhov odvodených od siete. Každá metóda bola vyhodnotená s rôznymi datasetmi od výskumnej organizácie. Taktiež rozšírili ich metódu pre vizuálnu analýzu siete emailových vírusov.

Ďalej Chapanond, Krishnamoorthy, Yener v [3] použil sietové metriky a spektrálnu analýzu k analýze či už orientovaného alebo neorientovaného grafu emailov, ktorý skonštruoval zmenou prahovej hodnoty (napr. počtom vymenených emailov medzi užívateľmi). Ich výskum je postavený na vytvorení emailového grafu a štúdiu jeho vlatnosťí či už pomocou teórie grafov alebo technikami spektrálnej analýzy. Grafová teoretická analýza zahŕňa výpočet niekolkých grafových metrík, ako napríklad rozdelenie podľa stupňov, priemerný pomer vzdialenosťí, zhlukovací koeficient alebo kompaktnosť emailového grafu. Hodnoty metrík v dátovej sade emailov spoločnosti Enron porovnali aj s inými emailovými dátami.

Jednou z univerzálniejsích prác je aj práca autorov Guanting Tang, Jian Pei, and Wo-Shun Luk [4]. Je to stručný prehľad hlavných výskumných snáh o analýzu mailov a popis metód, ktoré sa pri tejto analýze používajú. Nie len čo sa týka analytických alebo implemetnačných úloh, ale aj nástrojov, ktoré nám pri analýze vedia pomôcť. Aby zdôraznili rozdiely medzi analýzou mailov a bežnou analýzou textu, organizujú prieskum do piatich ľažších úloh a to: detekcia nevyžiadanej pošty, kategorizácia emailov, analýza kontaktov, analýza vlastností emailovej siete

a vizualizácia emailov. Tieto úlohy sú vlastne začlenené do rôznych spôsobov používania emailov. Systemaicky preskúmavajú bežne používané techniky a tiež budujú diskusiu o dostupných softwarových nástrojoch.

Na rozdiel od ostatných prác, Afra Abnar, Mansoureh Takaffoli, Reihaneh Rabbany, Osmar R. Zaiane [5] definovali vlastnú metodiku pre analýzu sociálnej siete a definovali *Structural social role mining framework*, ktorý je navrhnutý pre identifikáciu štrukturálnych rolí, pre identifikáciu zmien v sieti a analýzu dopadu zmien na sieť. Definujú základné sociálne roly v sieti a navrhujú metodológie pre ich identifikáciu. Pre identifikáciu týchto rolí využívajú klasické prostriedky analýzy sociálnych sietí a tiež navrhujú nové metriky zahrňujúc napríklad Betweenness centrality založenú na komunitách. Z tejto práce som vychádzala pri pomenovaní rolí zo siete a implementovala techniky pre ich identifikáciu.

Ďalšou prácou, ktorou som sa inšpirovala bola práca autorov Kudělka, Horák, Zehnalová [6], ktorá prezentuje analytický nástroj, ktorý bol vytvorený pre analýzu hlbších vzťahov v emailových dátach. Tieto vzťahy zahrňujú vzťahy založené na interakcii viacerých užívateľov v tíme. Analytické metódy popísané v práci sú založené na dvoch faktoroch. Prvým faktorom je kontext, čo je skupina viacerých užívateľov v kombinácii so slovami použitými v komunikácii. Druhým faktorom je časový interval, v ktorom bola začatá komunikácia. Práca prezentuje metódy pre väženie komunikácií, užívateľov a vzťahov, ako aj metód pre hľadanie komunít asociovaných so špecifickým kontextom.

### 3 Emailová komunikácia

#### 3.1 Stručná história emailu

Za počiatky emailovej komunikácie možno považovať priližne rok 1965, kedy bola správa prenášaná medzi sálovými počítačmi pracujúcich v režime zdieľania času na univerzite *Massachusetts Institute of Technology*.

Od tejto doby preša emailová komunikácia značným vývojom. Emaily, tak ako ich poznáme dnes, sú definované štandardom špecifikácie RFC2822 a sú prenášané pomocou komunikačných protokolov.

#### 3.2 Štruktúra emailu

Každý email sa skladá z dvoch častí - z tzv. hlavičky (*header*) a tela emailu (*body*).

Hlavička emailu je generovaná automaticky pri vytvorení emailu a sú do nej postupne vkladané informácie zo serverov, cez ktoré správa prechádza (tzv. MTA). Pre bežných užívateľov sú z hlavičky najdôležitejšie tieto údaje: predmet správy, čas odoslania, emailová adresa odosielateľa a prijímateľa. Ostatné údaje emailoví klienti (označovaní tiež ako MUA<sup>1</sup>) väčšinou nezobrazujú.

Pri vytváraní emailu emailovým klientom sú väčšinou do hlavičky vložené tieto záhlavia:

- **Date** - aktuálny čas počítača, ktorý vložil záhlavie
- **From** - adresa odosielateľa
- **Cc** - špecifikuje ďalších adresátov
- **Bcc** - umožňuje rozosielanie správy medzi viacerých adresátov
- **Priority** - priorita emailu, interpretácia sa lísi vzhľadom k MUA
- **Reply-To** - špecifikuje adresu, na ktorú je zaslaná prípadná odpoved
- **Subjekt** - predmet správy daný užívateľom
- **To** - udáva adresu príjemcu správy
- **Message-ID** unikátny identifikátor, ktorý je priradený MTA

Telo emailu obsahuje samotné dátá určené pre adresáta.

---

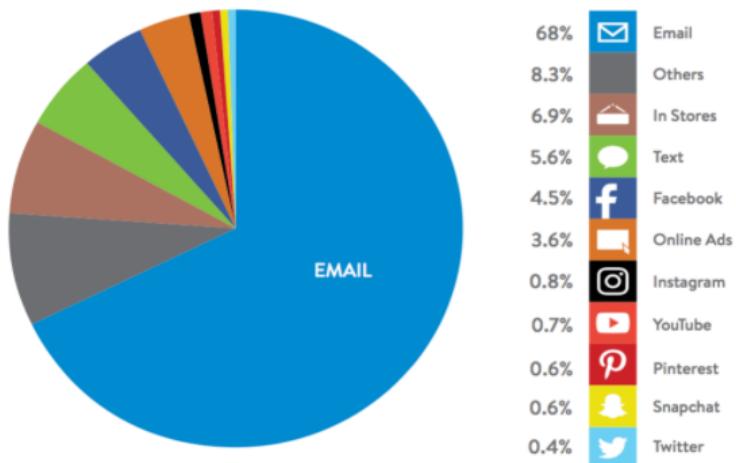
<sup>1</sup>MUA - Mail User Agent, program, ktorý používa užívateľ na rozosielanie a prijímanie emailov (napr. Outlook), tento program komunikuje s MTA (Mail Transfer Agent), ktorý sa stará o prenos emailov v prostredí verejnej siete Internet.

### 3.3 Emaily v súčasnosti

Emaily teda existujú už niečo cez 50 rokov, ich popularita je však stále veľká vďaka ich efektivite, extrémne nízkym nákladom a kompatibilite s množstvom typov zariadení. Ako jedna z najrozšírenejších typov komunikácie v dnešnej dobe, emaily sú široko rozšírené v každodennom živote. Napríklad, spolupracovníci diskutujú prácu cez emails, priatelia zdielajú sociálne aktivity a skúsenosti aj cez emails alebo veľké spoločnosti distribuujú reklamy práve pomocou emailov.

Aj keď by mnohí tvrdili, že éra emailov už je dávno preč a sú stále viac nahradzane novými sociálnymi sieťami, nové výskumy ukazujú opak. Napríklad výskum z roku 2016 od spoločnosti Bluecore [7] ukazuje, že email je stále populárny aj u mladších generácií, hlavne na formálnu komunikáciu.

V tomto výskume boli spotrebiteľia pýtaní, akú formu komunikácie preferujú pri komunikácii so značkami (interntovými obchodmi, na firemnú komunikáciu a celkovo formálnu komunikáciu). Prevažná časť opýtaných si vybrała email (68%).



Obr. 1: Akú formu komunikácie preferujete na formálnu komunikáciu?

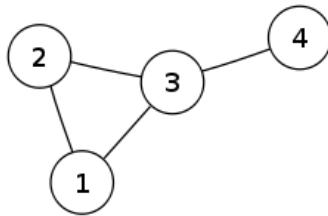
## 4 Definície a klasifikácie

V tejto kapitole popisujem všetky teoretické pojmy a metódy, ktoré v tejto práci spomínam a používam. V tejto kapitole budem používať matematické názvy podľa kontextu, v ktorom sa budem nachádzať.

### 4.1 Graf

- **Neorientovaný graf**

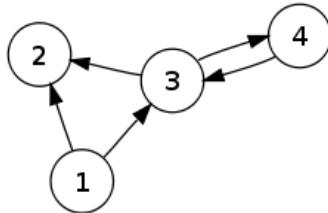
Neorientovaným grafov rozumieme usporiadanie dvojicu  $G = (V, E)$ , kde  $V$  je neprázdna množina *vrcholov* a  $E$  je neprázdná množina *hrán* - množina (niektorých) dvojprvkových podmnožín množiny  $V$ .



Obr. 2: Neorientovaný graf

- **Orientovaný graf**

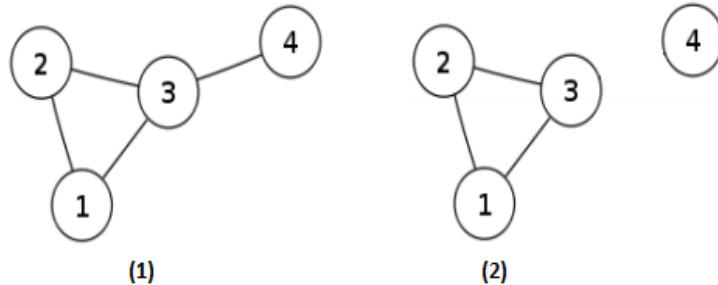
Orientovaným grafov rozumieme usporiadanie dvojicu  $G = (V, E)$ , kde  $V$  je množina *vrcholov* a množina orientovaných *hrán* je  $E \subseteq V \times V$ .



Obr. 3: Orientovaný graf

## 4.2 Súvislosť grafu

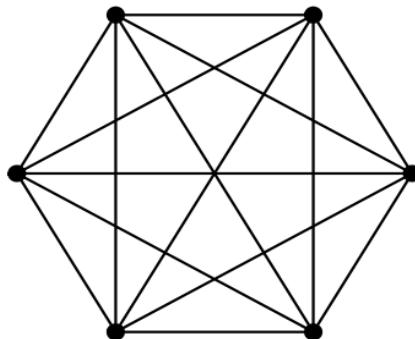
Hovoríme, že vrchol  $v$  je *dosiahnuteľný* z vrcholu  $u$ , ak v grafe existuje sled z vrcholu  $u$  do vrcholu  $v$ . Graf nazveme *súvislý*, ak pre každé dva vrcholy  $u, v$  je vrchol  $v$  dosiahnuteľný z vrcholu  $u$ . V opačnom prípade je graf *nesúvislý*.



Obr. 4: Súvislý (1) a nesúvislý graf (2)

## 4.3 Úplný graf

Úplný graf na  $n$  vrcholoch je neorientovaný graf, ktorý má hranu medzi každými dvoma vrcholmi. Počet jeho hrán je  $m = n(n - 1)/2$ .



Obr. 5: Úplný graf

## 4.4 Stupeň vrcholu

Stupeň vrcholu je počet vrcholov spojených s týmto vrcholom hranou, inými slovami: počet jeho susedov. V orientovanom grafe sa ešte rozlišuje vstupný a výstupný stupeň vrcholu podľa toho, koľko hrán z vrcholu vychádza alebo do neho vchádza.

## 4.5 Cesta

Cesta je postupnosť vrcholov v grafe taká, že medzi každými dvoma vrcholmi cesty je hrana a vrcholy sa neopakujú. V orientovaných grafoch sa ešte rozlišuje smer cesty, pričom orientácia hrán je stále rovnaká. Dĺžka cesty je počet hrán, ktoré obsahuje.

## 4.6 Uzavretá cesta

Uzavrená cesta, kružnica v neorientovanom a cyklus v orientovanom grafe, je cesta, ktorá začína a končí v rovnakom uzle.

## 4.7 Komponenta grafu

Komponenta grafu je súvislá časť grafu a medzi vrcholmi z rôznych komponent neexistuje žiadna hrana.

## 4.8 Metriky

V tejto časti popisujem metriky, ktoré v rámci identifikácie rolí v sieti používam. Ďalšie informácie o metrikách, ich praktickom využití a ich ďalších variantách sú zhrnuté v kapitole 6.2.

### 4.8.1 Closeness centrality (Centralita blízkosti)

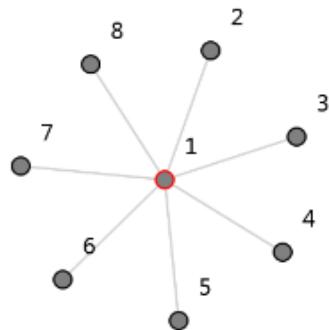
Táto centralita meria dôležitosť vrcholu grafu podľa priemernej hodnoty vzdialenosť od všetkých ostatných vrcholov v sieti. Aby dôležité vcholy mali vyššie číslo, je táto centralita počítaná ako inverzná hodnota tohto priemeru. Vrchol dôležitý podľa tejto metriky môže mať dobrý prístup k informáciám o ostatných vrcholoch alebo naopak môže ostatné vrcholy rýchlosťou ovplyvňovať.

Priemernú vzdialenosť vrcholu  $x_i$  od ostatných vcholov možno formálne zapísat ako  $l_i = \frac{1}{n} \sum_j d_{ij}$ , kde  $n$  je počet vrcholov v grafe a  $d_{ij}$  je najkratšia cesta medzi vrcholmi  $x_i$  a  $x_j$ . Centralita je potom  $C_i = \frac{1}{l_i}$ .

### 4.8.2 Betweenness centrality (Centralita medziľahlosti)

Táto centralita sa odlišuje od ostatných uvedených. Jej hodnota pre vrchol je počet najkratších ciest medzi každými dvoma vrcholmi v grafe, na ktorých hodnotený vrchol leží. Pokiaľ medzi vrcholmi v sieti tečú nejaké informácie alebo sa posielajú správy, hodnota tejto metriky vyjadzuje, aké množstvo informácií cez daný vrchol prejde. Táto centralita je tiež názorný príklad toho, že každá metrika počíta dôležitosť vrcholu úplne inak. Vrchol s vysokou centralitou medziľahlosti môže mať malý stupeň a nemusí ležať blízko ostatných vrcholov, stačí, keď cez neho prechádza veľa najkratších ciest. To môže nastať, pokiaľ vrchol je most medzi dvoma alebo viacerými komponentami v grafe, v extrémnom prípade pokiaľ je v strede grafu v tvare hviezdy (viď obrázok).

Betweeness vrcholu  $x_i$  možno spočítať ako  $B_i = \sum_{st} \frac{n_{st}^i}{g_{st}}$ , kde  $g_{st}$  je počet všetkých najkratších ciest medzi vrcholmi  $x_i$  a  $x_j$  a  $n_{st}^i$  je počet najkratších ciest, ktoré naviac vedú cez vrchol  $x_i$ .



Obr. 6: Graf v tvare hviezdy

#### 4.8.3 Modularita

Modularita je metrika, ktorá udáva rozdiel medzi počtom existujúcich hrán medzi vrcholmi rovnakého typu a počtom takých hrán v náhodne vytvorenom grafe v pomere ku všetkým existujúcim hranám. Vrcholy rovnakého typu sú tie, ktoré patria alebo majú patriť do rovnakej skupiny alebo triedy (komunity).

$$Q = \frac{1}{2m} \sum \sum_{i,j} \left[ A_{ij} - \frac{d_i d_j}{2m} \right] \delta(c_i, c_j)$$

Def: Modularita

## 5 Sociálna siet

Sociálna siet je množina sociálnych subjektov (uzly siete, spravidla jednotlivci alebo organizácie), ktoré sú prepojené jedným, alebo viacerými špecifickými druhmi vzájomnej závislosti, ako sú príbuzenstvo, priateľstvo, vzájomnosť, vízie, odpor, konflikt, obchod a pod. Sociálna siet z pohľadu teórie grafov je definovaná ako graf  $G(V, E)$ , kde  $V$  je množina entít (uzlov) a  $E$  je množina vzťahov (hrán) medzi týmito entitami.

Entity grafu môžu byť rôzne (základníci, jednotlivci, webové stránky, bankové účty, creditné karty, produkty). Nie je pravidlom, že len sociálna siet ako ju pozná mnoho ľudí je sociálnou sieťou aj formálne. Prvky sociálnej siete môžu ma napríklad aj skupina spolupracujúcich ľudí.

### 5.1 História sociálnych sietí

Pod pojmom sociálna siet si väčšina ľudí v dnešnej dobe predstaví služby ako *Facebook*, *Twitter* a pod. Tento pojem ale vznikol dlho pred vznikom internetu a dnešných sociálnych sietí. Prívlastok sociálny, ktorý sa v dnešnej dobe často vyniecha, je dôsledkom pôvodu analýzy sociálnych sietí. V druhej polovici 20. storočia sa simultánne v rôznych oblastiach skúmania vzťahov a chovania objavil nový pohľad na vzťahy medzi sociálnymi jednotkami a to ako na siet, graf. Preto prví predstavitelia analýzy sociálnych sietí boli pôvodne sociológovia alebo psychológovia (napríklad Moreno, Cartwright, Newcomb, Bavelas) a antropológovia (Barnes, Mitchell). Prvé použitie termínu "sociálna siet" sa pripisuje Barnesovi (1954).

V 30. rokoch 20. storočia psychiater Moreno rozvíjal sociometriu, predchodcu dnešnej analýzy sociálnych sietí. Vypytoval sa ľudí na priateľské vzťahy a skúmal, ako tieto vzťahy ovplyvňujú ich chovanie. Potom vynášiel (sám to tvrdil) tzv. *sociogram*, čo je diagram reprezentujúci ľudí ako body a vzťahy medzi ľuďmi ako úsečky, teda dnešnú sociálnu siet. Tento pojem sa ale začal používať až neskôr. Pomocou neho hľadal výrazné a izolované osoby v spoločnosti.

Zhruba o 20 rokoch neskôr antropológ Barnes začal skúmať, ako ovplyvnia vzťahy medzi ľuďmi nielen jednotlivcov, ale aj spoločnosť ako celok a zameral sa na štúdium skupín, komunít. Na práci Barnesa a jeho spolupracovníkov naviazala na Univerzite na Harvarde skupina vedená Harrisom Whitom. Tá začala budovať matematickú teóriu okolo dôležitejších pojmov zo sociálnych vied a umožnila tieto javy matematicky vyjadriť, merať a modelovať.

V druhej polovici 20. storočia sa rozšírilo povedomie o sociálnych sieťach a metódy sa začali používať aj v ďalších oboroch ako ekonómia, biológia, doprava atd.

### 5.2 Analýza sociálnych sietí

Analýza sociálnych sietí je interdisciplinárna veda s koreňmi v sociológii, psychológii, štatistike a teórie grafov. Analýza sociálnej siete chápe sociálnu siet ako systém prepojenia uzlov (individuálnych aktérov) prostredníctvom hrán (ich vzťahov). Možno teda povedať, že nadvázuje na matematickú teóriu grafov a metódy sietovej analýzy. Výsledkom analýzy môže teda byť mapa

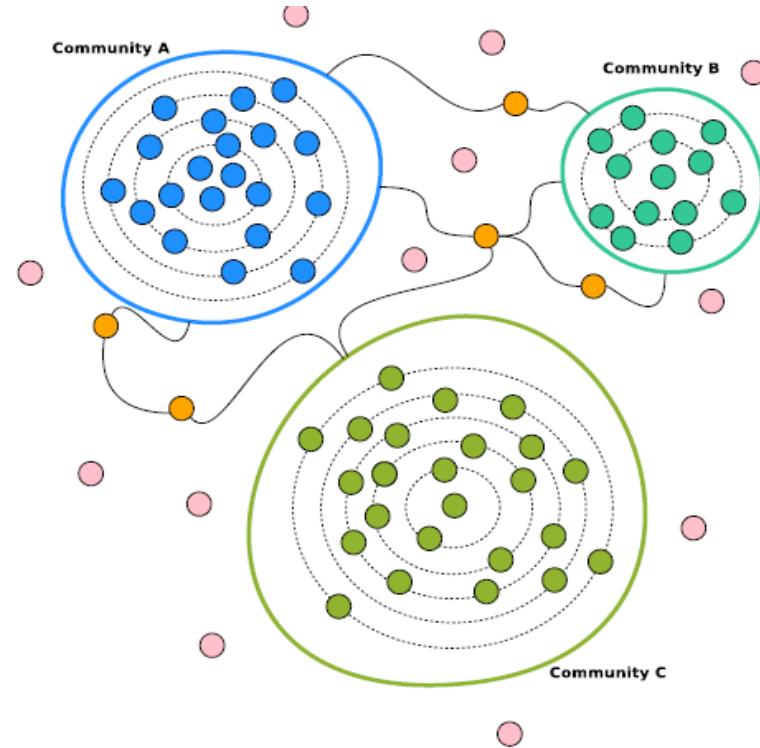
graficky znázorňujúca všetky prvky skúmaného sociálneho systému a ich vzťahy (resp. vybrané charakteristiky jednotlivých vzťahov vyjadrené vhodným spôsobom graficky). Charakteristikou môže byť napríklad vzájomná sympatia či antipatia alebo pravidelná vzájomná komunikácia alebo spolupráca.

Analýza sociálnych sietí vystupuje napríklad ako základná technika v rámci modernej sociológie, antropológie, sociálnej lingvistiky, geografie, sociálnej psychológie, ekonómie a biológie rovnako ako populárna téma pre výskum.

### 5.3 Komunity v sociálnych sieťach

Sociálne siete sú riedke grafy zložené z hustých podgrafov. Tieto husté podgrafové sú nazývané komunity. Najčastejšia definícia komunity: *Komunita je zhľuk uzlov, kde počet vnútorných hrán v komunite je väčší ako počet vnokajúcich hrán – mimo komunity.* [8]

Algoritmy pre dolovanie komunit sú založené na spojoch medzi uzlami, ktoré naznačujú spojenie dvoch entít. Napr. SCAN (Structural Clustering Algorithm for Networks) je metóda pre detekciu komunit v súvislosti na to, ako uzly zdieľajú svojich susedov len s ohľadom na priame spojenie. Teda ak sú dva uzly spojené a tiež zdieľajú rozumné množstvo ich susedov, patria do rovnakej komunity.



Obr. 7: Sieť s viacerými komunitami

## 5.4 Detekcia komunít

Detekcia komunít je proces identifikácie zhľukov uzlov siete silne prepojených medzi sebou a menej silne prepojených so zvyškom siete. Detekcia komunít v grafoch má za cieľ identifikovať moduly a ich prípadnú hierarchickú organizáciu.

Problém detektie komunít vyžaduje rozdelenie siete do komunít husto prepojených uzlov, pričom uzly patriace do odlišných komunít sú len slabo prepojené. Presné formulácie tohto optimalizačného problému sú známe ako výpočtovo neriešiteľné. Vyhľadávanie rýchlych algoritmov pritiaholo veľký záujem vďaka zvyšujúcej sa dostupnosti rozsiahlych sieťových dátových súborov a vplyvu sietí na každodenný život. Môžeme rozlišovať niekoľko typov algoritmov detektie komunít: *rozdelenacie* algoritmy - tie detekujú spojenie vnútri siete a postupne ich odstraňujú zo siete, *algomeratívne* algoritmy - zlúčujú podobné uzly a postupne komunity podľa spoločných črt a *optimalizačné* metódy sú postavené na maximalizácii objektívnej funkcie. Kvalita rozdielov vplývajúcich z týchto metód sa často meria takzvanou modularitou. Je to hodnota v intervale od -1 do 1, ktorá meria hustotu spojov vnútri komunít v porovnaní s prepojeniami medzi komunitami.

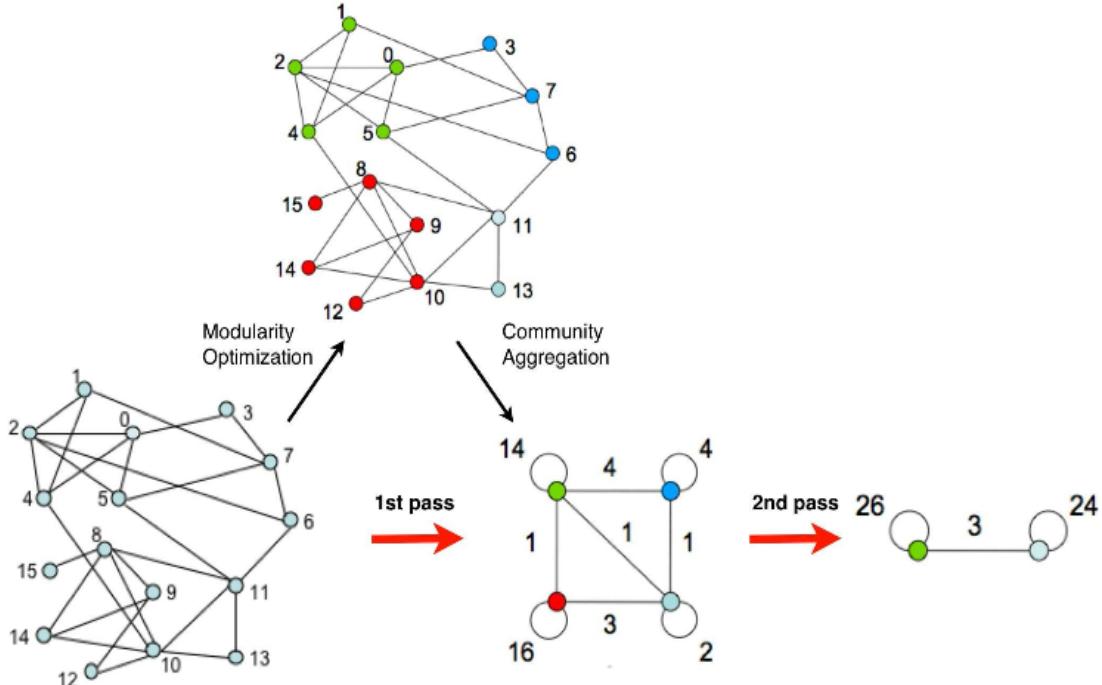
### 5.4.1 Louvainov algoritmus pre detekciu komunít

Veľmi oblúbeným a rýchlym algoritmom pre detekciu komunít je Louvainova metóda, ktorú navrhli Blondel, Guillaume, Lambiotte a Lefebvre [9]. Je to jednoduchá metóda pre extrakciu komunitnej štruktúry veľkých sietí. Je to heuristická metóda, ktorá je postavená na optimalizácii modularity. Je preukázané, že prekoná všetky ostatné známe metódy detektie komunít, pokiaľ ide o čas výpočtu. Navyše kvalita detekovaných komunít je veľmi dobrá.

Výpočet algoritmu je rozdelený do dvoch fáz, ktoré sa iteratívne opakujú. Predpokladajme, že začíname s váženou sieťou s  $N$  uzlami. Ako prvé označíme každý uzol siete inou komunitou. Takže v tomto prvotnom rozdelení je toľko komunít, ako je uzlov. Potom pre každý uzol  $i$  uvažujeme susedov  $j$  a vyhodnotíme prírastok modularity, ktorý by nastal, ak z sme odstránili uzol  $i$  z jeho komunity a priradili by sme ho do komunity uzla  $j$ . Uzol  $i$  je potom vložený do komunity, pre ktorú je tento prírastok najvyšší, ale len ak je tento prírastok kladný. Ak nie je možný žiadny kladný prírastok, uzol  $i$  ostáva vo svojej komuniti. Tento proces je aplikovaný opäťovne a sekvenčne pre všetky uzly kym sa nedosiahne žiadne zlepšenie a prvá fáza je kompletná. Prvá fáza končí, keď je dosiahnuté lokálne maximum modularity, keď žiadny uzol už nemôže zlepšiť modularitu. Je taktiež dôležité, že výstup algoritmu záleží na postupe, v ktorom sú uzly brané do úvahy. Výsledky algoritmu ale naznačujú, že usporiadanie uzlov nemá významný vplyv na získanú modularitu. Zoradenie však môže ovplyvniť výpočtový čas. Problém pri výbere objednávky preto stojí za to študovať, pretože by mohol poskytnúť dobrú heuristiku na zvýšenie výpočtového času.

Druhá fáza algoritmu spočíva vo vytvorení novej siete, ktorej uzly sú komunity nájdené počas prvej fázy algoritmu. K tomu, aby sa nová sieť vytvorila, vähy spojení medzi novými uzlami sú

dané sumou váh prepojení medzi uzlami korešpondujúcich dvoch komunít. Spojenia medzi uzlami tej istej komunity vedú k slučkám v novej sieti. Kedž je druhá fáza kompletnej, je možné znova aplikovať prvú fázu algoritmu na výslednú váženú siet a proces opakovať. Pri konštrukcii sa počet komunít znižuje pri každom priechode. Proces sa opakuje, kým nie sú žiadne ďalšie zmeny a dosiahne sa maximálna modularita.



Obr. 8: Vizualizácia krokov Louvainovho algoritmu.

Každý priechod je tvorený dvomi fázami: prvá, kde je modularita optimalizovaná tým, že umožňuje len miestne zmeny komunít a druhá, kde nájdené komunity sú agregované tak, aby bolo možné vytvoriť siet komunít. Priechody sú opakované iteratívne kým nie je možný žiadny nárast modularity.

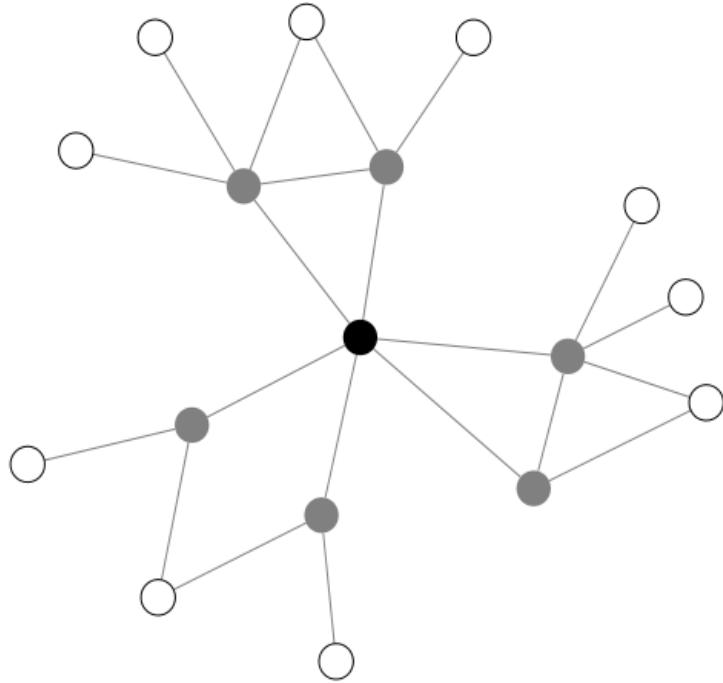
## 5.5 Ego siet

Ego siet je siet tvorená uzlami, ktoré sa nazývajú aj *alter* uzlami, ktoré sa formujú okolo určitého uzla, ktorý sa nazýva *ego*. Toto ego sa niekedy zo siete vynecháva za účelom analýzy zmien siete. To záleží od danej analýzy [10]. Ego je individuálny ústredný uzol. Siet môže mať toľko ég, kolko má uzlov. Egá môžu byť osoby, skupiny, organizácie alebo celé spoločnosti.

### 5.5.1 Konštrukcia ego siete

V tejto práci konštruujem ego siet tak, že k uzlu, ktorý bol vybraný ako ego, sa pridajú hrany tak, aby spájal doposiaľ neprepojené komponenty siete. To môže pomôcť k analýze toho, či by

sa daný človek hodil do vodcovskej pozície na základe jeho starých spojení. To, že sú k nemu pridané nové spojenia by bolo prirodzené, ak by sa daný jednotlivec dostal naozaj do vodcovskej (alebo inak egocentrickej) role. Príkladom môže byť napríklad to, keď sa v spoločnosti hľadá nový projektový manažér a z danej siete môžeme detektovať, či sa človek s danými vlohami nenachádza aj v aktuálnom tíme, ale na nižšej pozícii.



Obr. 9: Príklad ego siete.

## 6 Metódy analýzy sociálnych sietí

### 6.1 SSRM - Framework pre detekciu štrukturálnych rolí v sociálnych sieťach

Afra Abnar, Mansoureh Takaffoli, Reihaneh Rabbany, Osmar R. Zaiane [5] definovali *Structural social role mining framework*, ktorý je navrhnutý pre identifikáciu štrukturálnych rolí, pre identifikáciu zmien v sieti a analýzu dopadu zmien na sieť. Definujú základné sociálne roly v sieti(menovite Leader, Outermost, Mediator, Outsider).

#### 6.1.1 Rola v kontexte SSRM

*Sociálna rola* je sice základný sociologický pojem, ale stále neexistuje žiadny konsenzus v jej definícii. Podľa SSRM je rola je považovaná za pozíciu jednotlivca v spoločnosti. Informácie o sociálnej sieti sú kategorizované do štrukturálnych a neštrukturálnych vlastností. Štrukturálne vlastnosti sú príbuzné ku konštrukcii grafu ako sú spojenia entít (hrany), štruktúra susedov a pozícia entity v tejto štruktúre. Ale neštrukturálne vlastnosti sú ostatné informácie, ktoré neodrážajú konštrukciu grafu ako atribúty entít a spojení. SSRM definuje rolu v sieti ako: Rola entity v sieti je to, ako sa entita správa voči ostatným a jej vplyv na atribúty a štruktúry ostatných entít.

#### 6.1.2 Roly definované v SSRM

Ludské siete sú vnútorme zložené z viacerých komunít. V sociálnej sieti s viacerými komunitami, vlastnosti uzlov kolísu podľa toho, či je existencia komunít dostatočná alebo zanedbateľná. Z pohľadu sociálnej siete, uzol môže byť centrom celej siete, ale nie centrom v jeho komuniti. SSRM sa teda zameriava na štúdium sociálnych sietí s predpokladom existencie komunít v sieti, ako jej základnej črty.

V sociálnych sieťach môžu byť komunity explicitné alebo implicitné. Explicitné komunity sú postavené nezávisle na jej členoch a sú založené na množine pravidiel. V tomto prípade, ľudia sa stanú členmi tejto komunity častejšie až po zformovaní komunity. Zamestnanci firmy alebo študenti sú príkladom dvoch explicitných komunít. Zatiaľ čo formácia implicitných komunít tažko závisí na jej členoch a spojeniach. Tým pádom neexistuje žiadna externá podmienka na vybudovanie implicitnej komunity. Implicitné komunity sú postavené postupne ako sa ľudia spoločne stretávajú. Napríklad, skupina priateľov, v ktorej nie je žiadne pravidlo pre správanie sa jednotlivcov, je príklad implicitnej komunity. V oboch prípadoch explicitnej aj implicitnej komunity, by mali existovať aj špeciálne jednotlivci, ktorí tieto komunity manažujú a kontrolujú. Napríklad v školskej triede je to učiteľ alebo inštruktor. Pre firmu to je manažér vo vedení a pre skupinu priateľov je to zase človek, ktorého komunikačné schopnosti prinášajú ďalších členov alebo posilňujú vzťahy medzi tými stálymi. Títo dôležití jednotlivci sú ešte viac výrazný, keď je komunita obrovská.

Podľa toho SSRM framework definuje pre jednotlivcov v sociálnej sieti určité roly podľa ich vzťahov a pozícii v komunitách až po ich interakcie s ostatnými jednotlivcami. Z perspektívy komunít, v sieti existujú jednotlivci niekoľkých typov:

- so žiadnym vzťahom ku nejakej komunité
- so spojením s viacerými komunitami
- dôležitý členovia komunity
- bežný členovia komunity, ktorí formujú väčšinu
- nedôležitý členovia komunity, ktorí nemajú na komunitu pozorovateľný efekt

Na základe týchto poznatkov SSRM definuje štyri základné roly - **leader**, **mediator**, **outermost** a **outsider**.

#### 6.1.2.1 Leader

Sú mimoriadni jednotlivci v zmysle centrality alebo významu v každej komunité. V reálnom svete bývajú títo členovia siete veliteľmi, riaditeľmi, manažérmi, vládcami, prezidentami, autoritami, administrátormi atď.

#### 6.1.2.2 Outermost

Je to časť menej dôležitých jednotlivcov v každej komunité, ktorých vplyv a efekt na komunitu sú nižšie ako vplyv väčšiny členov komunity. Miesta, kde sa môže outermost v sieti nachádzať sú periféria alebo hranice grafu.

#### 6.1.2.3 Mediator

Sú to jednotlivci, ktorí zohrávajú dôležitú rolu v spojení komunít v medzi sebou. Fungujú ako mosty medzi odlišnými komunitami. Do tejto skupiny patria vyjednávači, sprostredkovatelia alebo aj rozbočovače v sieti.

#### 6.1.2.4 Outsider

Sú to jednotlivci, ktorí nie sú spojení so žiadnou komunitou v sieti. Bud majú takmer rovnaké prepojenie k rôznym komunitám alebo majú len veľmi slabé väzby na komunity.

### 6.2 Identifikácia štrukturálnych sociálnych rolí

Majúc sieť s komunitami explicitne známymi alebo extrahovanými nejakým dolovacím algoritmom, následne popisujem metodológie pre identifikovanie definovaných štrukturálnych rolí.

### 6.2.0.1 Outsider

Najviac priamočiarou rolou pre identifikáciu je outsider. Je to jednotlivec, ktorý v sieti nepatrí do žiadnej komunity. Identifikácia tejto roly je tak celkom priamočiara.

### 6.2.1 Leader

Leader je v každej komunite výnimočný centrálny člen. Pre identifikovanie takýchto uzlov SSRM využíva metriku *closeness centrality*.

#### 6.2.1.1 Closeness centrality (Centralita blízkosti)

V súvislom grafe closeness centrality uzlu je metrika centrality v sieti, vypočítaná ako súčet dĺžok najkratších cest medzi uzlom a všetkými ostatnými uzlami v grafe. Čiže čím viac je uzol centrálnejší, tým bližšie je k ostatným uzlom.

$$C_i = \frac{1}{l_i} = \frac{n}{\sum_j d_{ij}}$$

Def: Closeness centrality

Pre každý uzol sa stanoví hodnota closeness centrality. Hodnoty closeness centrality sú blízke notmálnemu rozdeleniu, v ktorom 95% populácie dát patrí do intervalu  $[\mu - 2 \cdot \sigma, \mu + 2 \cdot \sigma]$

Leadri ležia na hornom chvoste distribučnej funkcie, a teda horný interval použijeme pre identifikovanie leadrov. A teda uzly, ktoré majú väčšiu hodnotu closeness centrality ako krajná hodnota tohto intervalu, sú identifikovaní ako leadri.

### 6.2.2 Outermost

Podobne ako pri role *Leader* pre identifikovanie outermostov sa využíva metrika closeness centrality. Outermosti budú ležať však na spodnom chvoste distribučnej funkcie closeness centrality.

A tak teda uzly, ktoré majú hodnotu closeness centrality nižšiu ako  $[\mu - 2 \cdot \sigma]$ , sú outermosti.

### 6.2.3 Mediator

Rolu mediator zastávajú tí jednotlivci, ktorí spájajú viacero komunit a sú tzv. spojmy medzi komunitami.

Pre identifikáciu mediátorov sa definujú metriky založené na metrike betweeness centrality a to: *LBetweeness - LBC* a *CBetweeness - CBC* a ďalej metriky, ktoré vyjadrujú koľko rozdielnych komunit uzol spája: *DSCount* a *DSPair*.

#### 6.2.3.1 LBeweeness

**LPath** - Pred definíciou LBeweeness je potrebné definovať LPath a to nasledovne: *LPath* je množina všetkých najkratších cest medzi lídrami dvoch rozdielnych komunit.

$$LPath = l | startNode(l) \in leaderSet(c_i) \wedge endNode(l) \in leaderSet(c_j) \wedge c_i \neq c_j$$

Def: Lpath

**LBetweenss** centralita pre uzol  $v$  -  $LBX(v)$  je počet jedinečných LPath ktoré obsahujú  $v$ . Ak pre každú cestu  $p \in LPath$  definujeme  $I_l(p, v) = 1$  ak  $v$  leží na  $p$ , inak  $I_l(p, v) = 0$  potom:

$$LB(v) = \sum_{p \in LPath} I_l(p, v)$$

Def: LBetweenss

#### 6.2.3.2 CBetweenss

CBetweenss počíta počet najkratších ciest medzi rozdielnymi komunitami.  $s_p$  a  $e_p$  označujú štartovací a koncový uzol najkratšej cesty  $p$ . Taktiež  $c_v$  označuje komunitu, do ktorej uzol  $v$  patrí. Množina všetkých najkratších ciest, ktoré spájajú rozdielne komunity:  $CPaths = \{p | c_{s_p} \neq c_{e_p}\}$ . Taktiež definujeme  $I_p(p, v) = 1$  ak  $v$  leží na ceste  $p$  a 0 keď neleží.

$$CB(v) = \frac{1}{2} \sum_{p \in CPaths} I_p(p, v)$$

Def: CBetweenss

#### 6.2.3.3 Normalizovaná verzia CBetweenss

Pravdepodobnosť nájdenia viac viditeľných mediátorov vo väčších komunitách je väčšia v porovnaní s menšími komunitami. Táto situácia sa stáva, pretože vo väčších komunitách je pochopiteľne viac uzlov, čo vedie k viacerým najkratším cestám medzi nimi. Pre kompenzáciu tohto efektu je definovaná normalizovaná verzia  $CBC$ :

$$NBC(v) = \frac{1}{2} \sum_{p \in CPaths} \frac{I_p(p, v)}{\min(|c_{s_p}|, |c_{e_p}|)}$$

Def: Normalizovaná verzia CBetweenss

Navrhnuté metriky  $CBC$  a  $LBC$  sú nevyhnutné pre identifikovanie mediátorov, ale nie sú dostatočné. Napríklad pre sieť pozostávajúcu z desiatich komunít a dvoch mediátorov  $M_1$  a  $M_2$ , kde oba ležia na sto najkratších cestách medzi komunitami majú oba rovnaké hodnoty  $CBC$ . Kdežto  $M_1$  spája dve rozdielne komunity, kým  $M_2$  spája všetkých 10. Pri takomto scenárii  $M_2$  spája komunity viac globálne a mal by byť skôr posudzovaný ako mediátor ako  $M_1$ . A tak

*SSRM* definuje tzv. metriku **skóre rozmanitosti**, ktorá označuje rozdielne komunity, ktoré sú prepojené cez uzol.

#### 6.2.3.4 Skóre rozmanitosti

Táto metrika ukazuje koľko rozdielnych komunit je spojených cez špecifický uzol  $v$ . Túto metriku definujeme v dvoch variantach:

1. **DSCount** - je definovaný ako počet rozdielnych komunit, ktoré sú spojené daným uzlom.

Nech  $I_d(c_i, v) = 1$  ak  $\exists p \in CPaths : s_p \in c_i \wedge v \in p$ . Potom DSScount uzla  $v$  je definované ako:

$$DS_{count}(v) = \frac{1}{2} \sum_{c_i} I_d(c_i, v)$$

Def: DSCount

2. **DSPair** - Skóre rozmanitosti môže byť definované ako počet párov komunit, ktoré majú najmenej jednu najkratšiu cestu medzi ich členmi, ktoré prechádzajú uzlom  $v$ . Definujeme  $I_d(c_i, c_j, v) = 1$  ak  $\exists p \in CPaths : s_p \in c_i \wedge e_p \in c_j \wedge v \in p$

$$DS_{pair}(v) = \frac{1}{2} \sum_{c_i} \sum_{c_j \neq c_i} I_d(c_i, c_j, v)$$

Def: DSPair

Aj keď viac mediátorov môže mať rovnaké hodnoty jednotlivých metrík, môžu sa odlišovať napríklad v počte komunit, ktoré spájajú. *SSRM* to berie do úvahy a definuje tzv. *mediacy score* ako násobok normalizovanej CBetweeness a skóra rozmanitosti:

$$MS(v) = NCB(v) \cdot DS_{count}(v)$$

Def: Mediacy score

### 6.3 Brokerage roly

Jednoducho povedané, *brokerage* sa vyskytuje tam, keď jeden aktér siete poskytuje most medzi dvoma inými aktérmi, ktorí medzi sebou inak prepojení nie sú. Koncept *brokerage* rôl bol použitý vo veľmi veľa iných kontextoch, záleží len na jeho formalizácii. Aj keď je *brokerage* tradične konceptualizovaný ako dynamický fenomén, identifikácia *brokerage* rôl sa často využíva aj v oblasti statických spoločenských vzťahov.

Jedným známym kontextom pre *brokerage* je prípad obchodných vzťahov. V tomto prostredí, tito jednotlivci alebo organizácie, politické entity, ktorí boli schopní previezť tovar z jedného

miesta na druhé a kontrolovať ich rozšírenie, zohrávali kľúčovú rolu v udržiavaní obchodu na regionálnej a kontinentálnej úrovni. S prostredkováním kontaktov medzi vzdialé tretie strany (ktoré si nemôžu vmeniť informácie inak), títo aktéri povolili uvoľnenie kritických, priestorovo lokalizovaných zdrojov naprieč roziahlym územím, čo usnadňovalo rast zložitejších spoločností. Kým *brokerage* vo výmenných sieťach má dôležité systematické následky, jeho efekt na individuálnej úrovni bol oceňovaný viac intenzívne sociálmi (napr. v [11] [12] [13]).

Je zrejmé, že *brokerage* sa môže vyskytnúť v mnohých nastaveniach a povaha *brokerage* procesu samotného sa líši od kontextu. Vširšom zmysle tento proces spadá pod tri triedy - *transfer brokerage*, v ktorom *broker (ego)* vede informácie a iné zdroje od jedného jednotlivca k druhému, ktorí nie sú priamo prepojení. Potom *matchmaking brokerage*, v ktorom ego predstavuje alebo inak umožňuje spojenie jedného jednotlivca k druhému a nakoniec *coordination brokerage*, v ktorom ego usmerňuje kroky ostatných a tak vyriešia svoje závislosti bez toho, aby museli byť priamo prepojení.

*Brokerage* je stav alebo situácia, v ktorej účastník spája inak neprepojených účastníkov alebo zapĺňa medzery alebo diery v sieti. [11] Na obrázku je *broker* alebo aj *sprostredkovateľ* zastúpený čiernym uzlom, ktorý vyplňuje dieru v sieti alebo spojuje ostatných jednotlivcov reprezentovaných bielymi uzlami, ktoré predtým neboli navzájom prepojené priamo.



Obr. 10: Príklad brokerage procesu

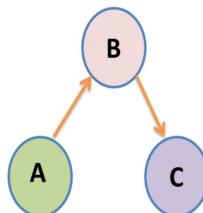
*Broker* môže prepojiť oddelené oblasti siete sociálnymi, ekonomickými alebo politickými aspektami a preto je jediný, kto má prístup k cenénym informáciám a zdrojom z rôznych oblastí siete. *Brokerage* je mechanizmus, ktorý umožňuje izolovaným či neprepojeným členom siete zdieľať informácie a zdroje a ekonomicky, politicky či spoločensky ovplyvňovať. [14]

Práve kvôli spojeniu a kontrole nad jedinečnými informáciami a zdrojmi medzi neprepojenými účastníkmi siete má aktér, ktorý zohráva rolu sprostredkovateľa (*broker*) v sieti väčší prístup k informáciám a zdrojom v porovnaní s tými, ktorí sprostredkovateľmi nie sú. *Broker* (sprostredkovateľ) môže fažiť z tejto kontroly nad informáciami a zdrojmi, môže sa stat silnejší v sieti a môže vykazovať zvýšenú efektivitu vo svojej práci. [14]

Detailnejšiu kategorizáciu *Brokerage* rôl predstavili Gould a Fernandez [11], kde predstavili koncept *brokerage* typológie. Táto typológia delí *brokerage* do piatich typov na základe smeru toku informácií - tokov v sieti - a rozdeľuje aktérov do vzájomne sa vylučujúcich skupín, tried alebo organizácií. Typy sprostredkovateľov sú *liaison*, *itinerant*, *coordinator*, *gatekeeper* a *representative*.

### 6.3.1 Liaison

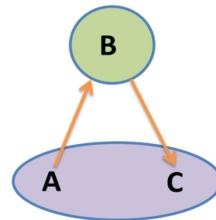
*Liaison brokerage* je *broker* spojenie medzi dvoma rozdielnymi skupinami, do ktorých on nepatrí. Na obrázku je *broker* (B) spojený s dvoma skupinami (A a C), ale nie je súčasťou ani jednej tejto skupiny a teda tvorí spojenie medzi aktérom A a aktérom C.



Obr. 11: Liaison brokerage

### 6.3.2 Itinerant

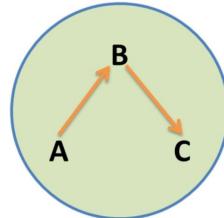
Pri tomto type *brokerage* roly dvoja neprepojení aktéri (A a C) patria do jednej skupiny, kým *broker* (B) patrí do inej skupiny. *Itinerant brokerage* je tiež nazývaný *consultant brokerage*, pretože broker sa chová ako konzultант pre oboch nespojených aktérov tej istej skupiny.



Obr. 12: Itinerant brokerage

### 6.3.3 Coordinator

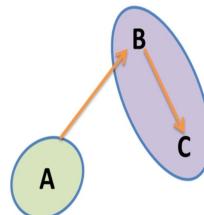
V role *coordinator* všetci traja aktéri patria do rovnakej skupiny a sprostredkovanie informácií a zdrojov sa deje v rámci skupiny.



Obr. 13: Coordinator brokerage

#### 6.3.4 Gatekeeper

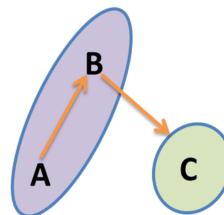
V tomto typie *brokerage* roly *broker*(B) a jeden z dvoch neprepojených aktérov (C) patria do jednej skupiny, kým iný neprepojený aktér (A) patrí do rozdielnej skupiny. *Broker* tohto typu kontroluje prichádzajúce informácie a zdroje v rámci jeho skupiny a robí rozhodnutia a tom, či majú alebo nemajú neprepojení aktéri v skupine prístup k informáciám a zdrojom.



Obr. 14: Gatekeeper brokerage

#### 6.3.5 Representative

*Representative* rola je podobná role *gatekeeper* role, *broker* (B) a jeden nespojený aktér (A) patra do jednej skupiny kým ten druhý nespojený aktér (C) patrí do inej rozdielnej skupiny, ale smer toku informácií alebo zdrojov je rozdielny.



Obr. 15: Representative brokerage

#### 6.3.6 Identifikácia brokerage rolí

Päť typov *brokerage* rôl reprezentujú unikátne sociálne roly zapuzdrujúce elementárny aspekt aktérovej štrukturálnej pozície v danej sieti. Jeden jednotlivec však môže zohrávať viac *brokerage* rolí naraz. Preto Gould & Fernandez [15] kvantifikovali celkovú participáciu jednotlivca v *brokerage* rolách pomocou *brokerage* skóra. Formálne definovali *brokerage* v grafe reprezentujúcim asymetrickú reláciu  $R$ : Nech  $a$  je *broker* medzi  $b$  a  $c$  iba ak  $bRa$ ,  $aRc$  a  $a\bar{R}c$ , kde  $bRa$  indikuje, že  $b$  je prepojené s  $a$  reláciou  $R$  a  $b\bar{R}c$  je negácia  $bRc$ . S touto definíciou, *brokerage* skóre sa vypočíta súčtom počtu kôľko krát táto podmienka platí pre špecifickú kombináciu spojenia aktérov. To znamená, že ak nejaký aktér  $x$  zohráva pozíciu *coordinator* dva krát a pozíciu *representative* tri krát, tak aktér bude mať skóre pre pozíciu *coordinator* = 2, pre pozíciu *representative* = 3 a jeho celkové *brokerage* skóre bude 5.

Formalizácia *brokerage* rolí podľa Goula a Fernandeza je definovaná pre siete, v ktorých sú spojenia (hrany) orientované, čiže reprezentujúce vzťahy, pre ktoré môžeme rozlíšiť odosielateľa a prijímateľa. Kedže v mojej koncepcii siete, kde jednotlivé uzly sú členovia tímu a hrana medzi nimi je práve vtedy, keď medzi nimi prebehla konverzácia, moja vytvorená sieť je neorientovaná. Zovšeobecnenie na neorientovanú sieť je celkom jasné; s takýmito dátami, každá hrana je považovaná za obojsmernú. Aj keď toto prináša jednu dôležitú zmenu originálnej formalizácie: v prípade neorientovaných vzťahov nemôžeme rozlíšiť rolu *gatekeeper* od roly *representative*, pretože neprítomnosť obojsmerných vzťahov redukuje tieto dve roly do jednej a *brokerage* skôr bude pre tieto dve roly identické. [15]

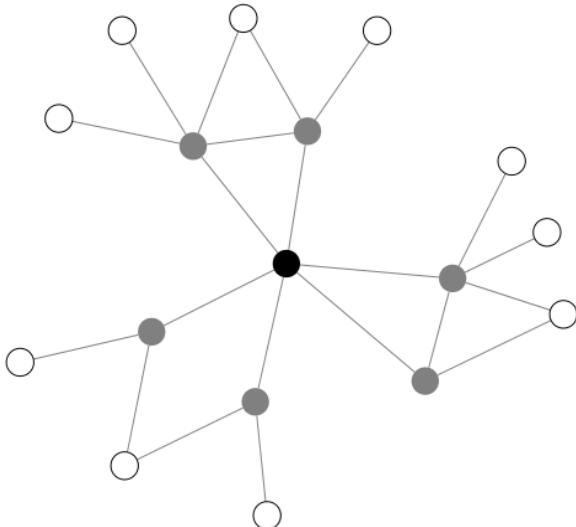
## 6.4 Analýza ega

Analýza ego sietí sa stáva stále viac dôležitou s rastom sietí. Je oveľa jednoduchšie v obrovských sietach analyzovať ego a jeho okolie ako celú sieť ako celok. Napríklad ak jeden človek má priemerne 5 blízkych osôb, potom v meste s populáciou desať tisíc ľudí bude päťdesiat tisíc priateľských väzieb. A ak by sme chceli študovať známosti? Riešením by bol výber podmnožiny obyvateľov mesta a ich *alter* uzlov.

Ďalšou odpovedou na otázku, prečo študovať ego siete, je to, že niekedy nás nezaujíma sieť ako celok alebo komunity a podobne, ale zaujímajú nás dôležitosť alebo inak zaujímaní jednotlivci (lídri, umelci, tínedžeri a pod.) Siet ega je zaujímatá, pretože je zdrojom informácií, sociálnej podpory, prístupu ku zdrojom, vplyvu a ďalších faktorov.

### 6.4.1 Veľkosť ego siete

Veľkosť ego siete je jednoduchá, ale veľavrvná charakteristika. Definuje ju stupeň ego uzla, alebo teda počet *alter* uzlov ega. Hovorí o sociálnej podpore, prístupu k informáciám a zdrojom.



Obr. 16: Veľkosť ega - stupeň uzla: 6

#### 6.4.2 Kompozícia ego siete

Čo sa týka kompozície ego siete, môžme sledovať podobnosť medzi egom a jeho *alter* uzlami. Pre reprezentáciu podobnosti sa používa *homofília*. Môžeme predpokladať, že existuje vzťah medzi nejakým javom a tým, či ego zdiela so svojimi *alter* uzlami nejakú vlastnosť (profesia, vzdelanie a pod.). Napríklad je prirodzené, keď niekto, kto zastáva pozíciu CFO (Chief Financial Officer) je obklopený ľuďmi, ktorí riešia finančie alebo napríklad politici bývajú obklopení členmi rovnakej politickej strany.

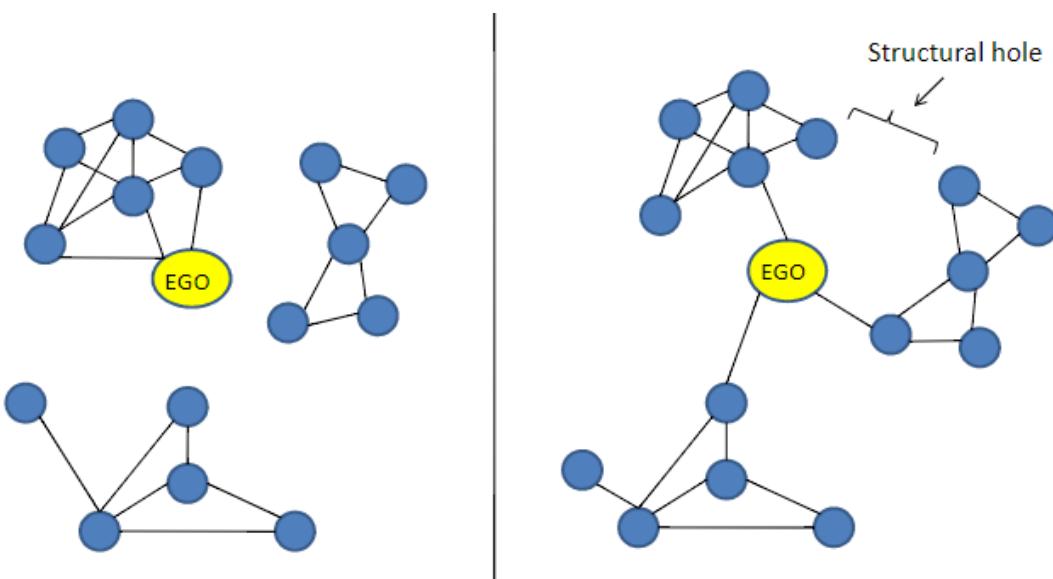
Pre identifikáciu homifílie som využila prítomnosť komunit v sieti a použila som *Krackhardt-Sternov E-I index* [16].

$$\frac{E-I}{E+I}$$

- $E$  je počet spojení s členmi inej skupiny (komunity)  $I$  je počet spojení s členmi rovnakej skupiny (komunity)
- nadobúda hodnoty od -1(homofília) do +1(heterofília)

#### 6.4.3 Štruktúra ego siete

Štrukturálna analýza sa opiera o informácie, či existujú alebo neexistujú spojenia medzi *alter* uzlami ego uzla. Princíp spočíva v tom, že nedostatok spojení medzi *alter* uzlami môže priniesť určité benefity samotnému egu. Tento princíp sa v analýze sociálnych sietí nazýva princíp štrukturálnych dier (ang. *structural holes*). Medzi benefity, ktoré prinášajú štrukturálne diery egu patria prístup k novým informáciám, moci alebo k slobode.



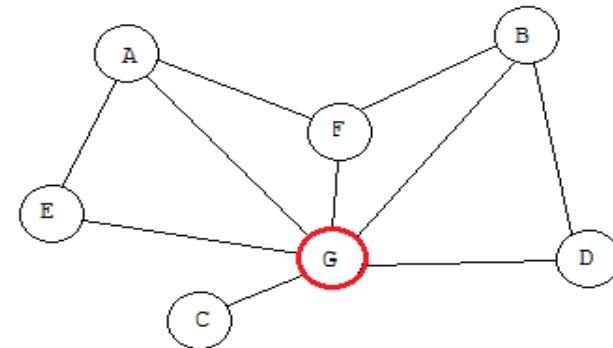
Obr. 17: Málo štrukturálnych dier vs. veľa štrukturálnych dier.

Koncept štrukturálnych dier je koncept analýzy sociálnych sietí vyvinutý R. S. Burtom. Predstavil tento pojem v snahe vysvetliť vznik rozdielov v sociálnom kapitále. Burtova teória

naznačuje, že jednotlivci majú isté výhody alebo nevýhody podľa toho, ako sú zakotvené v spoľočenských štruktúrach. Štrukturálna diera je chápaná ako medzera medzi dvoma jednotlivcami (chýbajúca hrana medzi uzlami), ktorí majú doplňujúce zdroje informácií. [17]

#### 6.4.3.1 Efektívna veľkosť

Burt predstavil mieru redundancie siete, Borgatti vyvinul zjednodušenú verziu efektívnej veľkosti pre nevážené siete [18]. Redundancia =  $\frac{2t}{n}$ , kde  $t$  je počet všetkých spojení v egocentrickej sieti (s výnimkou spojení k egu) a  $n$  je počet všetkých uzlov v egocentrickej sieti (s výnimoku ega). Táto formula môže byť modifikovaná pre výpočet efektívnej sily ego siete. Efektívna sila ego siete =  $n - \frac{2t}{n}$  (rozdiel počtu alter uzlov ega a sumy ich redundancií =  $6 - 1.33 = 4.67$ )



Uzol G je ego	A	B	C	D	E	F	Celkom
Redundancia	3/6	2/6	0/6	1/6	1/6	1/6	1.33

Obr. 18: Príklad výpočtu efektívnej sily

Čím viac je každý uzol odpojený od ostatných primárnych kontaktov, tým vyššia bude efektívna veľkosť. Tento indikátor nadobúda hodnoty od 1 (sieť poskytuje len jediné spojenie (hranu)) až do celkového počtu spojení, kedy každý kontakt (alter) je neredundantný.

## **7 Aplikácia**

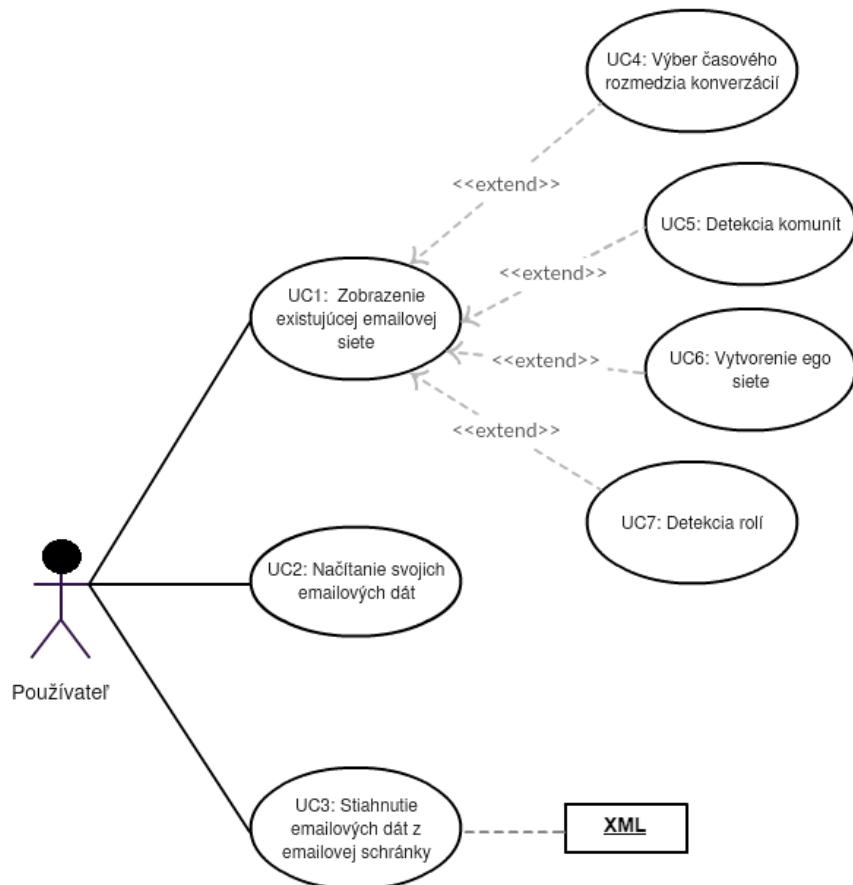
Táto kapitola obsahuje všetky podrobnosti o vývoji aplikácie, návrhu a ďalej špecifikáciách požiadavkov. Sú tu uvedené informácie o implementácii, návrhu, návrhových vzoroch, ale aj konštrukcii siete, predpríprave dát. Táto časť taktiež obsahuje diagramy najdôležitejších tried aplikácie alebo diagramy prípadov použitia.

### **7.1 Špecifikácia**

Aplikácia slúži ako užívateľské rozhranie na analýzu emailovej komunikácie a vizualizáciu analytických výstupov. Aplikácia umožňuje exportovať dátá z emailovej schránky alebo importovať vlastný XML súbor s emailovými dátami a ďalej s týmito dátami pracovať a zobrazovať siet emailovej komunikácie. Umožňuje vytvorenie ego-siete alebo detektovať vo vytvorennej siete komunity. Najdôležitejšou časťou aplikácie je detekcia štrukturálnych rolí v sieti, čiže detekcia dôležitých a nedôležitých členov emailovej komunikácie.

#### **7.1.1 Funkčné požiadavky**

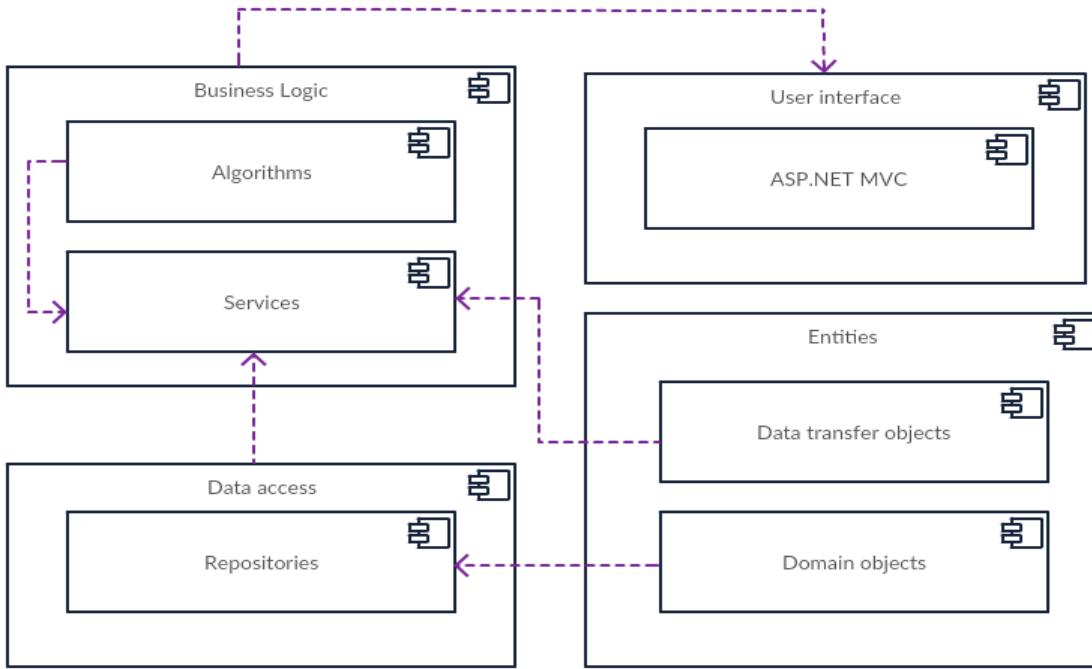
- Export dát z emailovej schránky
- Import vlastného XML súboru s emailovými dátami
- Zobrazenie informácií o emailovej sieti
- Vizualizácia emailovej siete
- Vytvorenie ego-siete
- Detekcia komunít
- Detekcia štrukturálnych rolí v sieti
- Výber časového rozmedzia emailových konverzácií



Obr. 19: UseCase Diagram

## 7.2 Návrh

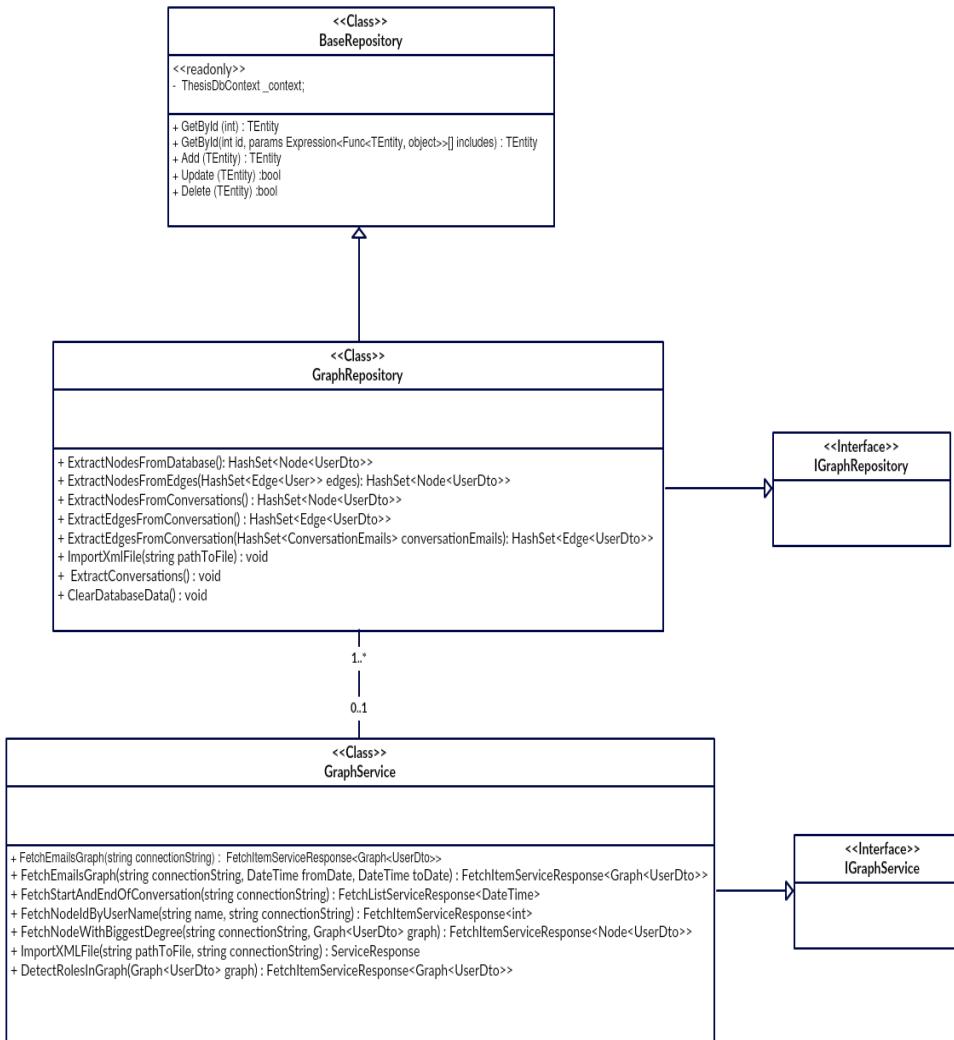
Aplikácia je vytvorená ako .NET aplikácia (veria .NET Frameworku 4.6). Je vytvorená ako trojvrstvová, pre uloženie dát sa používa SQL databáza. Najnižšia vrstva aplikácie slúži na získavanie dát z databázy, pre prepojenie s databázou a posielanie dát z aplikácie do databázy používam Entity Framework a používam tu návrhový vzor Repository. Od tejto časti je oddelená časť s business logikou a na najvyššej časti, ktorá slúži len na zobrazenie dát a komunikáciu s užívateľom, používam známy prístup Model-View-Controller.



Obr. 20: Diagram komponent znázorňujúci jednotlivé komponenty architektúry aplikácie

### 7.2.1 Návrhové vzory

**Repository** Návrhový vzor Repository je základným kameňom doménou riadeného návrhu. Model aplikácie teda nemá poňatie o tom, akým spôsobom je perzistovaný. O to sa stará práve Repository. Naviac práve vďaka tomu, že sa o persistenciu stará cudzí objekt, stačí poznať len jeho rozhranie a v prípad potreby ho ľahko nahradíť iným. [19]

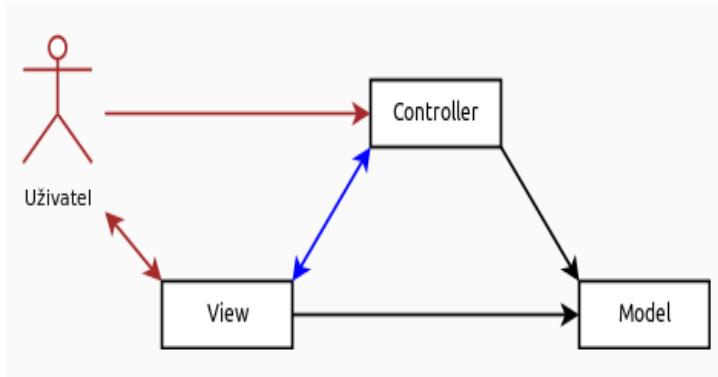


Obr. 21: Triedny diagram - Repository pattern

**Model-View-Controller** V aplikácii je použitý tradičný vzor Model View Controller (MVC). Je to jeden z najpoužívanejších a najobecnejších architektonických vzorov.

MVC rozdeľuje program do troch hlavných častí:

- **Model** - dátá a súvisiace operácie
- **View** - prezentácia dát (užívateľské rozhranie), obsahuje priamy odkaz na model, aby mohol jeho dátá prezentovať vonkajšiemu svetu
- **Controller** - riadi tok udalostí v programe, konkrétnie v tejto aplikácii kontrolery obsahujú len volanie metód z inej vrstvy aplikácie



Obr. 22: Model-View-Controller

### 7.3 Dôležité rozhodnutia

Pri navrhovaní aplikácie bolo potrebné urobiť niekolko dôležitých rozhodnutí.

#### 7.3.1 Dostupnosť dát

Pôvodne sa zvažovalo použitie aplikácie a analýzy dát nad verejne dostupnou anonymizovanou emailovou sadou. Emailových dát je ale veľmi málo a chcela som, aby sa výsledky práce dali overiť nie len inými analytickými nástrojmi, ale aj empiricky. Takže som využila to, že pracujem a moja emailová schránka teda nie je chudobná na maily. Navyše mi radi pomohli aj moji kolegovia a poskytli mi svoje emailové dáta. Tako som zozbierala reálne emailové dáta štyroch ľudí, o ktorých je známe ich postavenie v týme alebo aj dátum nástupu do práce. Takže výsledky daných algoritmov som vedela porovnať s reálnou situáciou v kolektíve.

#### 7.3.2 Webová vs. desktopová aplikácia

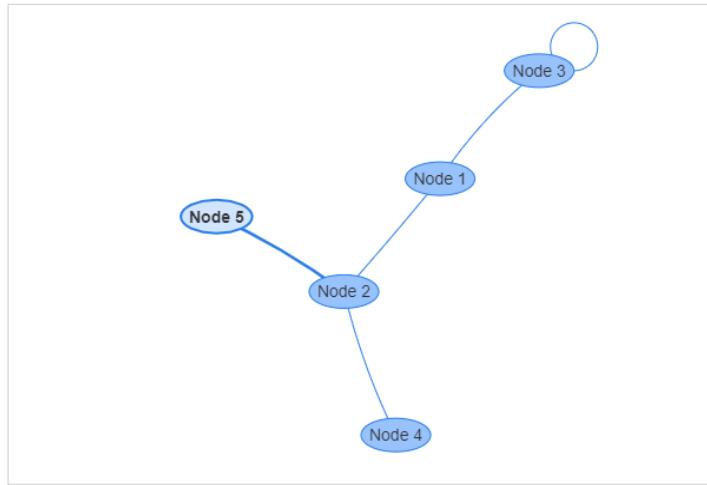
Bolo nutné sa rozhodnúť, či vyvíjať aplikáciu ako webovú alebo desktopovú. Ako platforma bola zvolená Microsoft .Net a programovací jazyk C#. Jednou z variant bola desktopová aplikácia vyvíjaná vo Windows Forms. WinForms je osvedčná technológia, je vyladená, v základe obsahuje veľké množstvo grafických prvkov. Toto sú výhody rozšírenej a dlho používanej technológií. Nevýhoda je ale práve zastaralosť a ťažkopádnosť v kreslení a spravovaní grafického rozhrania. Keďže ale doba ide dopredu a web a webové aplikácie sú stále viac používaniejšie a v súčasnosti existuje mnoho grafických knižníc pre vizualizáciu grafického rozhrania, rozhodla som sa aplikáciu vyvíjať ako webovú.

### 7.4 Použité knižnice

#### vis.js

Vis.js je dynamická vizualizačná knižnica. Knižnica je navrhnutá tak, aby bola ľahko ovládateľná a aby mohla spracovať veľké množstvo dynamických dát a umožňovala manipuláciu s dátami a

interakciu s nimi. Knižnica sa skladá z častí *DataSet*, *Timeline*, *Network*, *Graph2d* a *Graph3d*. Pre moju aplikáciu som používala len časť *Network*.



Obr. 23: Jednoduchá sieť vytvorená s použitím knižnice vis.js

```
<style type="text/css">
    #mynetwork {
        width: 600px;
        height: 400px;
        border: 1px solid lightgray;
    }
</style>

<script type="text/javascript">
    // create an array with nodes
    var nodes = new vis.DataSet([
        {id: 1, label: 'Node 1'},
        {id: 2, label: 'Node 2'},
        {id: 3, label: 'Node 3'},
        {id: 4, label: 'Node 4'},
        {id: 5, label: 'Node 5'}
    ]);

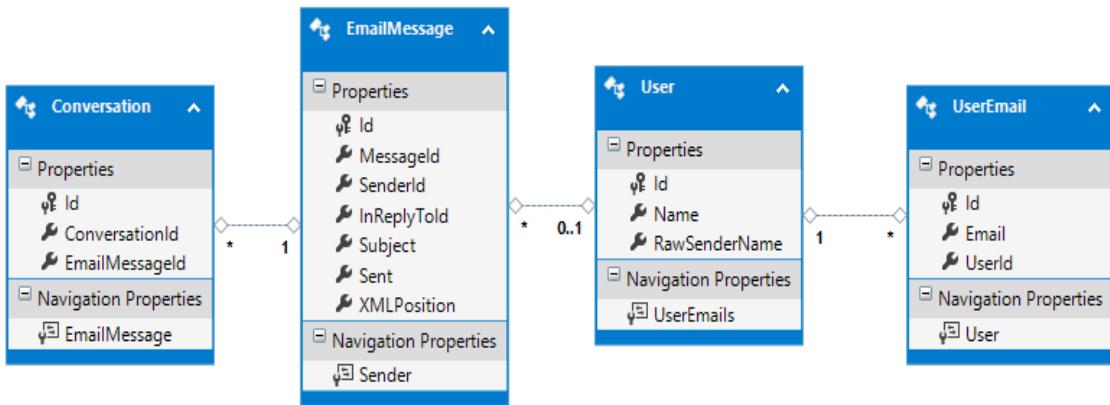
    // create an array with edges
    var edges = new vis.DataSet([
        {from: 1, to: 3},
        {from: 1, to: 2},
        {from: 2, to: 4},
        {from: 2, to: 5},
        {from: 3, to: 3}
    ]);

    // create a network
    var container = document.getElementById('mynetwork');
    var data = {
        nodes: nodes,
        edges: edges
    };
    var options = {};
    var network = new vis.Network(container, data, options);
</script>
```

Obr. 24: Príklad použitia knižnice vis.js

## 7.5 Import dát

Do aplikácie je možné nahrať XML súbor, ktorý je spracovaný uloženou SQL procedúrou, ktorá rozparsuje emailové dátá na jednodlivé entity - *User*, *EmailMessage*, *UserEmail* a *Conversation* a uloží ich do SQL databázy.



Obr. 25: Doménový model

## 7.6 Implementácia

Aplikácia je napísaná v jazyku C#, grafické rozhranie je naimplementované pomocou návrhového vzoru Model View Controller a graf bol vizualizovaný pomocou knižnice vis.js. Aplikácia bola vyvíjaná vo Visual Studiu 2017.

### 7.6.1 Metóda pre získanie emailových dát

Pre získanie emailových dát z emailovej schránky som naimplementovala metódu, ktorá sa pomocou protokolu IMAP pripojí na danú emailovú schránku a stiahne emails vo forme XML súboru.

EMAIL SERVER CONFIGURATION

Email		
veronika.uhrova@globallogic.com		
Username	veronika.uhrova@globallogic.com	
Password	*****	
Server address	Port	
imap.gmail.com	993	
<input checked="" type="checkbox"/> Use secure connection(SSL)		
<input type="button" value="Submit"/>		

Obr. 26: Príklad konfigurácie emailu pre získanie emailov

### 7.6.2 Konštrukcia siete

Rozdiel medzi príspom rôznych štúdií a mojím prístupom pri konštrukcii grafu z emailového datasetu je v konštrukcii komunikačnej siete. Ako základnú stavebnú jednotku siete som si zvolila **konverzáciu**. Inšpirovala som sa prácou autorov Kudělka, Horák, Zehnaloová [6]. Konverzácia je teda súbor emailov, ktorá začína jediným emailom, obsahuje najmenej 2 emaily a dvoch rôznych odosielateľov. Vrcholom siete (grafu) sa teda stane užívateľ, ktorý bol ako odosielateľ aspoň v jednej takejto konverzáции. Hrana medzi užívateľmi je zostrojená medzi užívateľmi, ktorí boli spolu v jednej konverzáции ako odosielatelia. Pre konverzáciu ešte ukladám čas jej začiatku, užívateľ si následne v aplikácii môže zvolať časový rozsah konverzácií.

### 7.6.3 Triedy pre graf, vrcholy a hrany

Pre uloženie siete v pamäti slúži generická trieda `Graph<T>`. Vrcholy a hrany drží ako `Dictionary<int, HashSet<Node<T>>`, čiže ako mapu vrcholov s ich susednými vrcholmi. Pre uloženie vrcholov a hrán grafu slúžia zoznamy hrán a vrcholov uložené ako `HashSet<Node<T>>` a `HashSet<Edge<T>>`. Trieda je generická preto, aby bolo možné vytvoriť graf pre rôzne entity. Triedy reprezentujúce vrcholy a hrany sú tiež generické a predstavujú ich `Node<T>` a `Edge<T>`. Pre identifikáciu komunity, bola vytvorená trieda `Community<T>`.

## 8 Experimenty

V tejto kapitole popisujem experimenty, ktoré som previedla nad rôznymi emailovými sadami v implementovanej aplikácii. Popisujem ako aj prípravu dát, tak aj import a ďalej vizualizáciu jednotlivých datasetov a tiež výsledky a porovnanie analytických metód. V prvej časti popisujem analýzu agregovanej sady ako celkovú analýzu tímu a v druhej časti popisujem analýzu jednotlivca.

### 8.1 Analýza tímu

Analýzu tímu som previedla na agregovanej emailovej sade, ktorá obsahuje email od štyroch jednotlivcov. Základné informácie o tomto datasete sú uvedené v tabuľke. Ďalej popisujem aj prípravu a import dát do aplikácie.

	počet
Emaily	21738
Konverzácie	4149
Požívateľia	370
Používateľia, ktorí sú aspoň v jednej konverzáции	273
Emaily poslané z pracovného emailu	20367

Tabuľka 1: Základné informácie o datasete

#### 8.1.1 Príprava a import dát

Ako už bolo spomenuté, pre samotnú analýzu som nepoužila žiadny verejne dostupný dataset, ale použila som svoju emailovú sadu a emailové sady od mojich troch kolegov z práce. V dobe, keď som od nich emaily požadovala, moja aplikácia ešte nebola hotová a tak sa pre získanie ich emailov použila externá aplikácia *TeamNet Data* [20]. Emailové sady boli poskytnuté v súbore formátu XML.

Kedže sme kolegovia a pracujeme spolu v tíme, nachádzali sme sa viac krát v rovnakej emailovej komunikácii, takže sa v sadách vyskytovali emaily duplikátne. Pre spracúvanie emailov som vytvorila SQL procedúru, ktorá jednotlivé XML súbor načíta, rozdelí emaily, používateľov a konverzácie na jednolivé entity a nainštalauje ich do SQL databázy. Pre odstránenie duplikátov som použila SQL skript, ktorý zaručil odstránenie duplikátnych položiek. S taktiež uloženými emailovými dátami používam pre import do aplikácie *Entity framework*, ktorý zaručuje prenos dát medzi aplikáciou a SQL databázou.

Po spustení aplikácie mám na výber stiahnuť si svoj XML súbor o svojej emailovej schránke, nainštalať ho do aplikácie a ďalej prezeráť svoje emailové dátá. Kedže emailové sady som už mala predpripravené, mohla som postupovať priamo k vizualizácii a analýze.

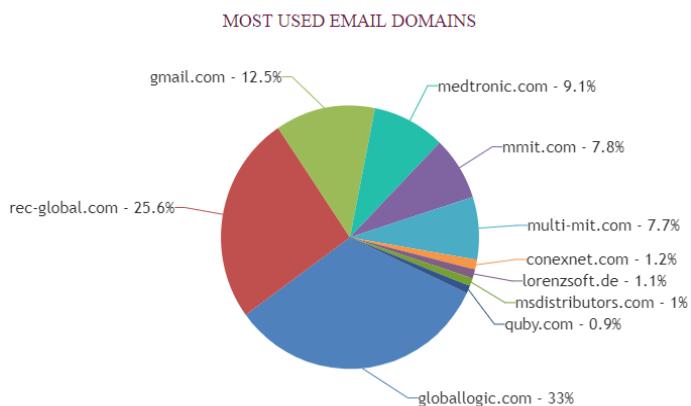
### 8.1.2 Vizualizácia datasetu

Aplikácia umožňuje výber vizualizácie jednotlivých členov tímu a tiež celého tímu celkovo. S výberom daného datasetu sa pre danú sieť zobrazia aj základné informácie ako počet emailov, odosielateľov emailov, najpoužívanejšie emailev domény a podobne.



Obr. 27: Základné informácie o tímovej sieti.

Na obr. 27 sú zobrazené základné informácie o agregovanej sieti. Štyria sme celkovo napísali 21738 emailov z toho bolo detekovaných 4149 konverzácií. Celkovo sa na emailovej komunikácii podieľalo 370 používateľov. Najviac emailov v agregovanej sade poslal používateľ *Tibor Palatka*. Hodina, kedy sa posielalo celkovo najviac emailov bola medzi ôsmou a deviatou hodinou ráno, čiže to je čas, kedy ľudia prídu do práce a prvé, čo urobia je, že si skontrolujú emaily.



Obr. 28: Najviac používané emailové domény.

Na obr. 28 sú zobrazené emailové domény, z ktorých používateelia najviac posielali emaily. Najviac emailov sa poslalo z domén *rec-global.com* a *globallogic.com*, ktoré sú oficiálne domény spoločnosti, v ktorej pracujem. Ostatné domény patria zákazníkom, s ktorými ako spoločnosť spolupracujeme.

Na obr. 29 je vizualizácia celkovej siete tímu v základnej podobe.



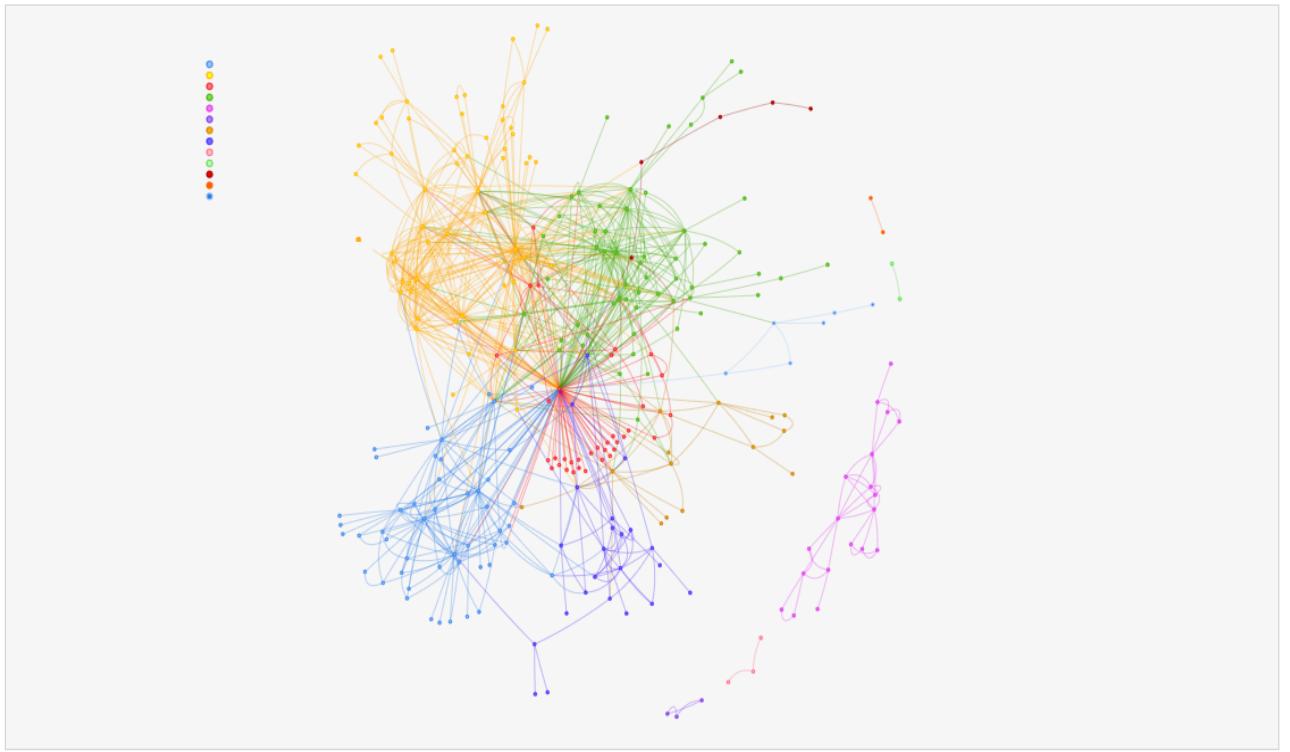
Obr. 29: Vizualizácia siete.

### 8.1.3 Detekcia komunít

Aplikácia umožňuje detektovať komunity v sieti a následne ich vizualizovať. Pre ďalšiu analýzu poskytujem v nasledujúcej tabuľke základné informácie o členoch tímu. Pre analýzu komunít využívam to, že každý nastúpil do práce v iný rok, takže na prelome týchto rokov by mal byť zaznamenaný nárast komunít.

Meno	Aktuálna pozícia	Dátum nástupu do firmy
Andrej Parimucha	Test engineer	1.7.2011
Veronika Uhrová	Junior Developer	1.9.2016
Andrej Matejčík	Medior Developer	1.6.2017
Tibor Palatka	Lead Developer	1.11.2010

Tabuľka 2: Informácie o členoch tímu



Obr. 30: Vizualizácia komunít v tímovej sieti za celkový čas

Celkovo bolo v sieti detekovaných 13 komunít. Najväčšia komunita má 57 uzlov, najmenšia má 2 uzly. Celkové zloženie komunít je zobrazené na nasledujúcom obrázku.

COMMUNITIES	
Community 1	48 nodes
Community 2	57 nodes
Community 3	35 nodes
Community 4	56 nodes
Community 5	19 nodes
Community 6	3 nodes
Community 7	14 nodes
Community 8	23 nodes
Community 9	3 nodes
Community 10	2 nodes
Community 11	5 nodes
Community 12	2 nodes
Community 13	6 nodes

Obr. 31: Rozloženie komunít v tímovej sieti za celkový čas

### 8.1.3.1 Zmeny komunít v čase

Zmeny komunít som zaznamenávala v dátumoch, kedy jednotliví členovia tímu nastupovali do práce. Prvým časovým úsekom, ktorý som sledovala bolo obdobie medzi 1.11.2010 - 1.7.2011, kedy bol zamestnancom spoločnosti zatiaľ len jeden z nás. Boli detekované dve komunity, rozloženie uzlov je zobrazené na obrázku 32



Obr. 32: Rozloženie komunít za prvý časový úsek

Ďalší interval, ktorý som zvolila bol interval medzi 1.11.2010 - 31.8.2016, ktorý reprezentuje čas, kedy boli zamestnancami dvaja jednotlivci. Detekovaných bolo 9 komunít.



Obr. 33: Rozloženie komunít za druhý časový úsek

Do príchodu ďalšieho kolegu (interval 1.11.2010 - 31.5.2017) bolo detekovaných 12 komunít, čiže mojim príchodom do firmy pribudli tri komunity, to znamená že príchodom ďalšieho kolegu počet komunít vzrástol o jednu komunitu. Na ďalšom obrázku je zobrazené rozloženie komunít v treťom intervale.



Obr. 34: Rozloženie komunít za tretí časový úsek

#### 8.1.4 Ego siet

Pre analýzu ego siete som vybrala troch jednotlivcov, z ktorých som vytvorila ego uzol a vytvorila pre nich ego siet. Porovnávala som jednotlivé metriky a počet komunít, ktoré spájali. Vyberala som uzly s rozdielnou centralitou, aby som mohla sledovať porovnanie medzi výsledkami. Porovnanie je možné vidieť v tabuľke.

Meno jednotlivca	Stupeň uzla	Počet prepojených komunít	Efektívna veľkosť	E-I index
Tibor Palatka	129	13	127	0.52
Andrej Parimucha	81	11	72	0.06
Attila Soltézs	20	11	18	0.4

Tabuľka 3: Informácie o vytvorennej ego sieti

Podľa veľkosti, čiže stupňa uzla môžem povedať, že uzol s veľkosťou 129 bude mať určie väčší sociálnu podporu od ostatných aktérov v sieti, väčší prísun k zdrojom a informáciám v porovnaní s uzlami s veľkosťou 81 alebo 20. Čo sa týka počtu prepojených komunít sú ale títo aktéri na tom podobne a spájajú porovnatelné množstvo komunít. Podľa E-I indexu môžem povedať, že ego siet každého aktéra je heterofílna a tak sa v svojej práci stretávajú so širokou skálou ľudí, nielen s ľuďmi z podobného okolia.

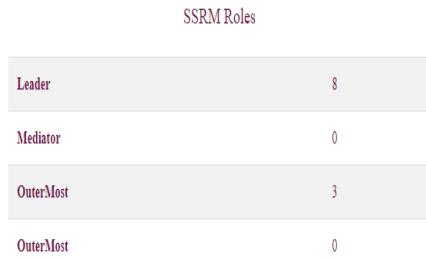
#### 8.1.5 Analýza rolí

Analýzu rolí som prevádzala na sieti s vytvoreným ego uzlom (ako ego uzol som zvolila uzol s najväčším stupňom).

##### 8.1.5.1 SSRM

V na obr. 35 sú zobrazené počty detekovaných štrukturálnych rôl v sieti. Uzly, ktoré boli de-

tekované ako *leader* zastávajú v reálnom živote pozíciu *Lead developer*, *Tester*, *Product owner*, *HR manager*, *Project manager*, *Sales manager*, *IT Consultant* a ďalší *Project manager*. Žiadny uzol neboli detekovaný ako *Mediator*. Tri uzly boli detekované ako *Outermost*, čo majú byť menej dôležitý jedinci a naozaj som sa s nimi ani s ich menom v práci nestretla.



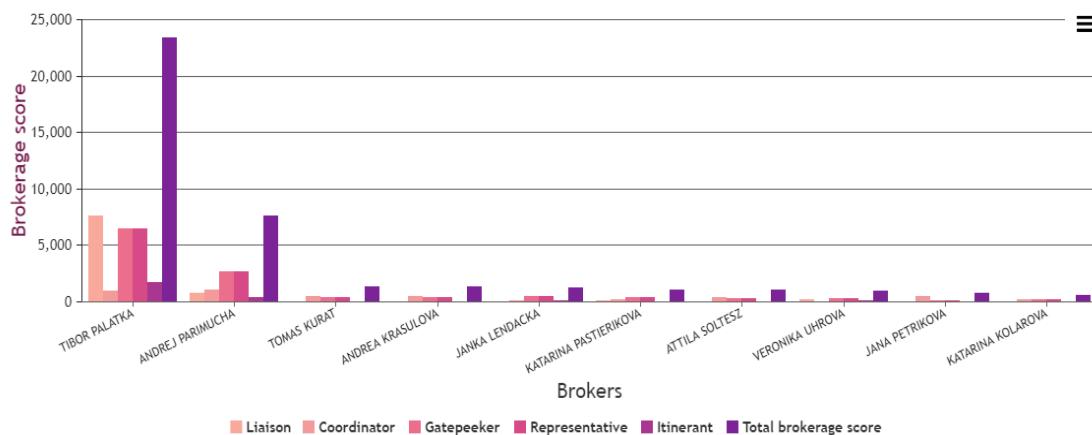
Obr. 35: Počet detekovaných štrukturálnych rôl

#### 8.1.5.2 Brokerage

V na obr. 36 a 38 je zobrazených desať najväčších *broker* aktérov v sieti spolu s grafom ich čiatkovým skóre pre každú *brokerage* rolu, ako aj celkové *brokerage* skóre.

TOP 10 BROKERS							
	Name	Coordinator	Itinerant	Gatekeeper	Representative	Liaison	Total
1	TIBOR PALATKA	976	1758	6486	6486	7636	23342
2	ANDREJ PARIMUCHA	1084	404	2708	2708	744	7648
3	TOMAS KURAT	498	12	420	420	26	1376
4	ANDREA KRASULOVA	460	14	408	408	40	1330
5	JANKA LENDACKA	116	130	480	480	0	1206
6	KATARINA PASTIERIKOVA	162	20	416	416	72	1086
7	ATTILA SOLTEZZ	344	6	326	326	40	1042
8	VERONIKA UHROVA	44	54	324	324	210	956
9	JANA PETRIKOVA	506	0	110	110	4	730
10	KATARINA KOLAROVA	204	6	176	176	20	582

Obr. 36: Desať aktérov s najväčším *brokerage* skóre



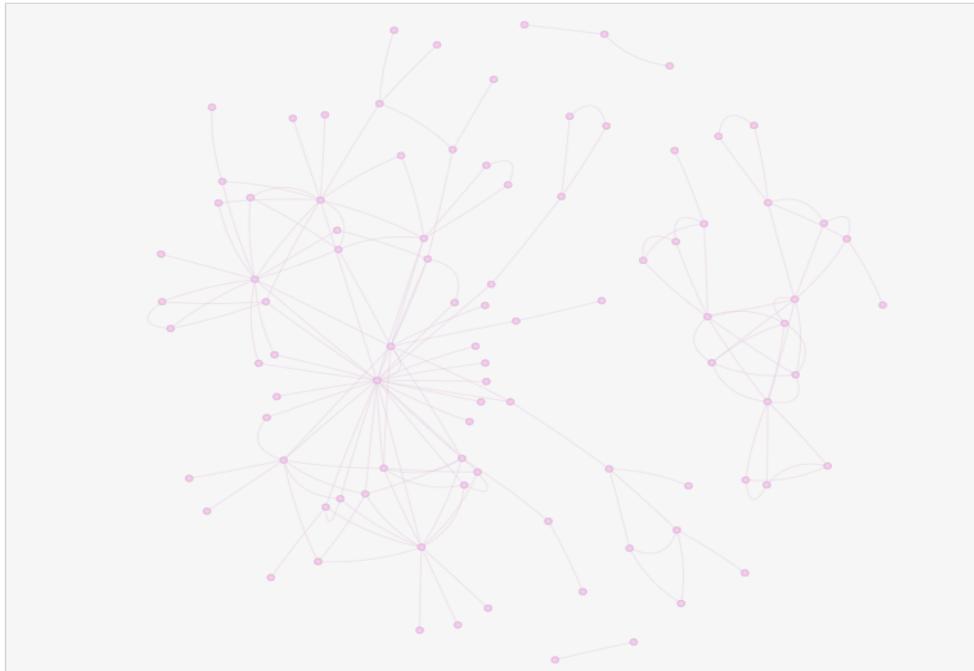
Obr. 37: Desaf aktérov s najväčším *brokerage* skóre - graf

## 8.2 Analýza jednotlivca

V tejto časti popisujem analýzu jednotlivca od získania emailových dát z emailového účtu, importu do aplikácie a zobrazenie výsledkov metód v aplikácii. Pre analýzu jednotlivca používam svoju emailovú sadu pre demonštrovanie získavania a importu emailov .

### 8.2.1 Príprava a import dát

Pre získanie emailových dát z emailovej schránky používam navrhnutú aplikáciu. Po zadaní emailového účtu, používateľského mena, hesla, adresu servera a portu získam XML súbor s emailami, ktorý naimportujem do aplikácie a zobrazí sa mi moja emailová sieť.



Obr. 38: Analýza jednotlivca - základná vizualizácia

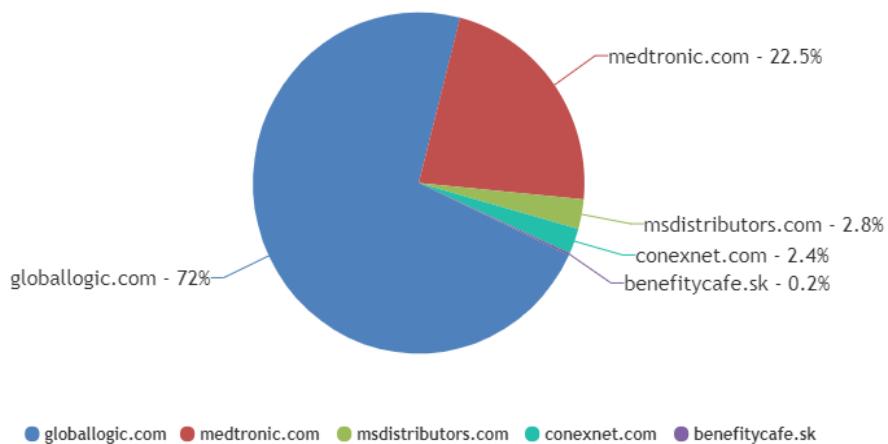
### 8.2.2 Informácie o datasete

Spolu s vizualizáciou datasetu sa zobrazia aj základné štatistiky o datasete. Moja emailová schránka obsahuje 3135 emailov, ktoré napísalo 155 používateľov a bolo detekovaných 512 konverzácií. Osoba, ktorá napísala najviac emailov bol používateľ *Attila Soltezs*, ktorý zastáva pozíciu *Project manager*. Hodina, kedy sa písalo najviac emailov bola medzi siedmou a ôsmou hodinou ráno. Celkový sumár spolu s najpoužívanejšími emailovými doménami je na obrázku 40.

## STATISTICS

NUMBER OF EMAILS	3135
NUMBER OF CONVERSATIONS	512
NUMBER OF USERS	155
BIGGEST EMAIL SENDER	ATTILA SOLTESZ
MOST EMAILS IN CONVERSATION	11
PEEK HOUR	7:00-8:00

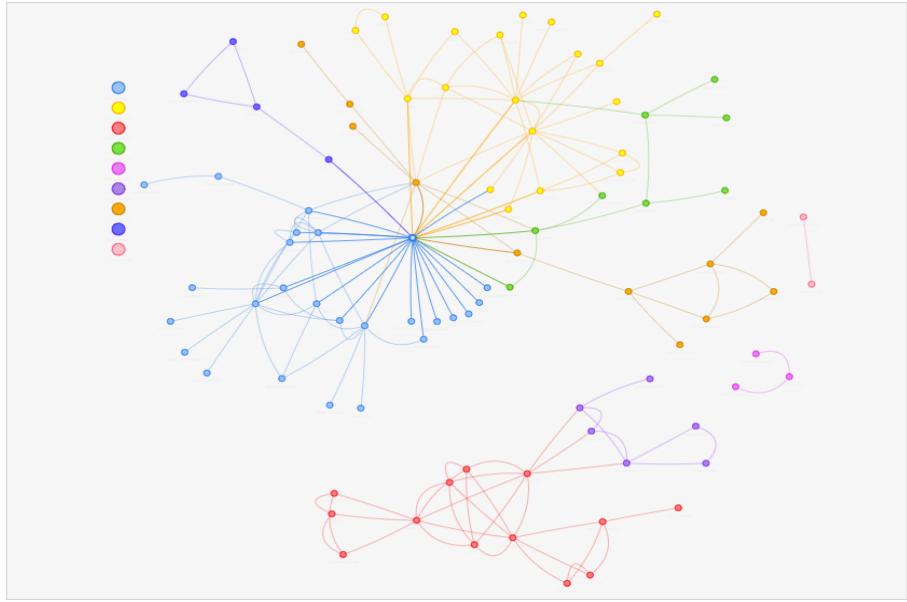
## MOST USED EMAIL DOMAINS



Obr. 39: Analýza jednotlivca - základné štatistiky

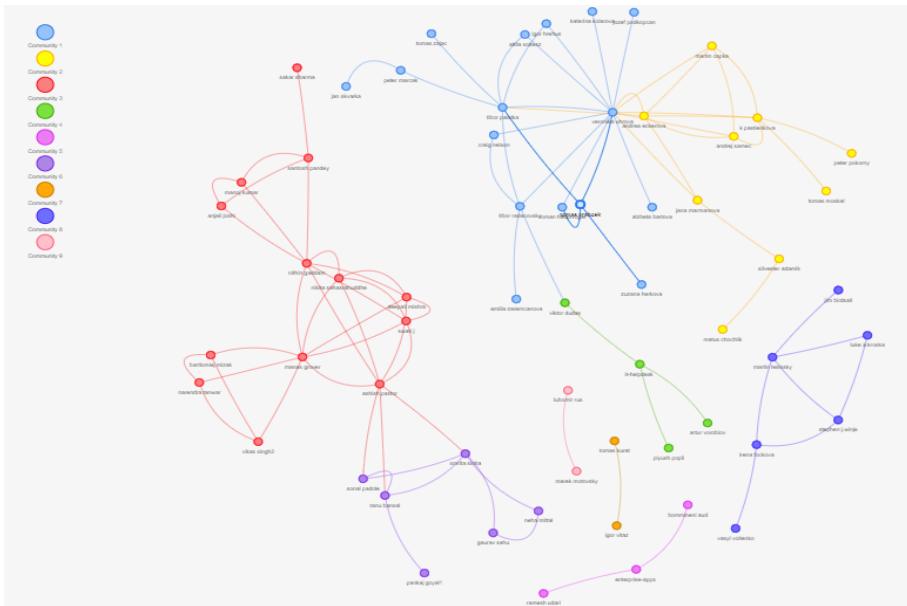
### 8.2.3 Detekcia komunit

Celkovo bolo v dataste detekovaných 9 komunít. Najväčšia komunita obsahuje 26 uzlov, najmenšia komunita obsahuje 2 uzly.

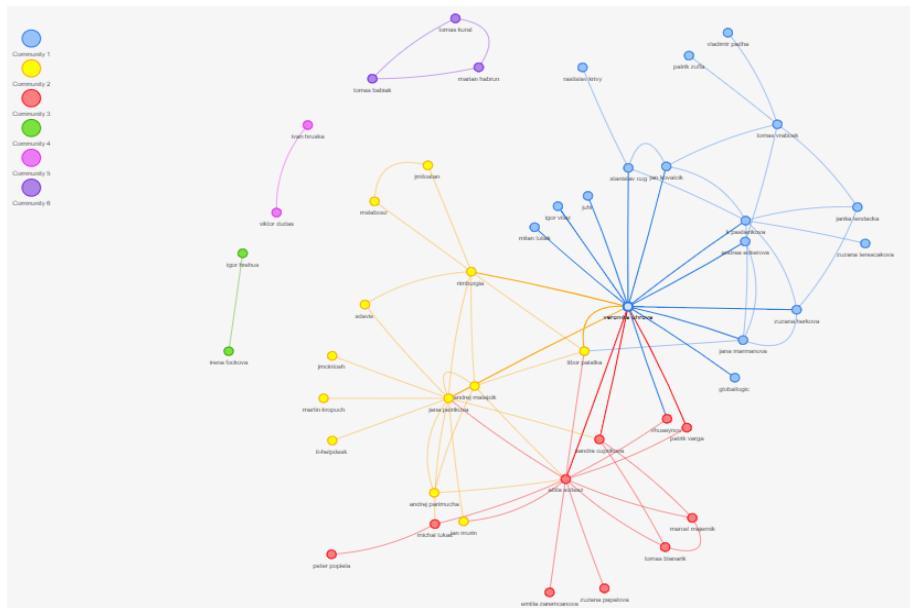


Obr. 40: Analýza jednotlivca - vizualizácia komunít

**8.2.3.1 Zmeny komunít v čase** Pre zmeny komunít v čase som využila obdobia, kedy som pracovala na iných projektoch a teda sa komunity môžu v týchto intervaloch lísiť. Obdobie prvého projektu je v intervale od 1.9.2016 do 1.3.2017 a obdobie druhého je od 1.3.2017 do súčasnej doby, povedzme do 31.3.2018.



Obr. 41: Analýza jednotlivca - Vizualizácia komunít v prvom časovom intervale



Obr. 42: Analýza jednotlivca - vizualizácia komunit v druhom časovom intervale

Ako vidieť na obrázkoch 41 a 42 zloženie komunit sa lísi, či už zložením jednotlivých uzlov, tak aj veľkosťou a počtom, čo je samozrejme prirodzené, keďže som komunikovala v týchto intervaloch s inými osobami. V prvom intervale bolo detekovaných 9 komunit, najväčšia komunita mala 16 uzlov, najmenšia 2 uzly. Čo sa týka druhého intervalu, bolo detekovaných 6 komunit, najväčšia komunita obsahovala 17 uzlov najmenšia rovnako 2 uzly.

#### 8.2.4 Ego sieť

Pre analýzu ega som vybrała cielene troch jednotlivcov zo siete, ktorí pracujú na rozdielnej pozícii. Vybrala som jedného developera, jedného projektového manažéra a testera. Z každého z nich vytváram v sieti ego a v nasledujúcej tabuľke porovnávam ich silu v sieti.

Aktér	Pozícia	Prepojené komunity	Efektívna veľkosť	E-I index
Veronika Uhrová	developer	8	28	-0.06
Attila Soltézs	projektový manažér	11	15	-0.25
Andrej Primucha	tester	4	5	0

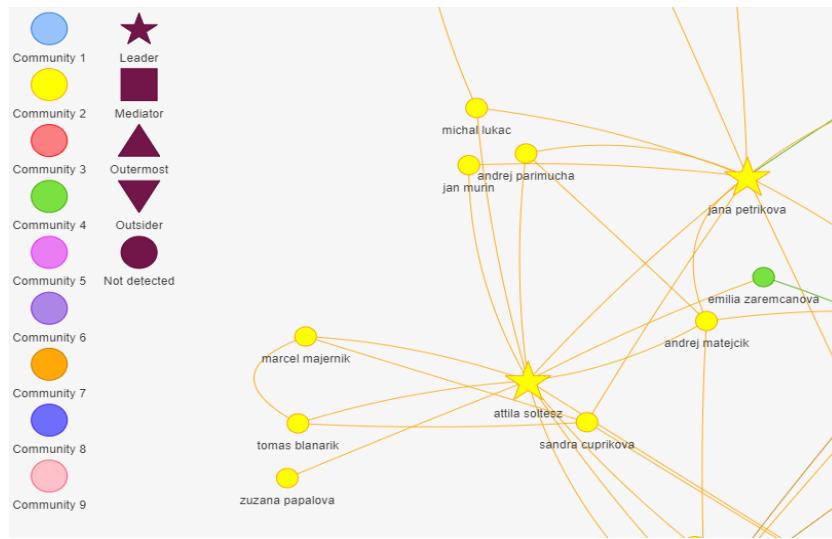
Tabuľka 4: Informácie o vytvorenej ego sieti

#### 8.2.5 Analýza rolí

Analýzu rolí som prevádzala na sieti s detekovaným ego uzlom, ktorý som zvolila podľa najväčšieho stupňa uzla.

### 8.2.5.1 SSRM

Analýza štrukturálnych rolí dopadla podľa očakávaní. Boli detektované tri uzly s roľou *leader* a jeden uzol s roľou *mediator*. Uzly detektované ako *leader* pracujú na pozíciach *Project Manager*, *Product owner* a *Developer*. Možno trochu nezvyčajné je, že mňa ako developera môj algoritmus identifikoval ako *leader* rolu. Keď to premietnem do reálneho života, kde je v tíme jediná žena, môžem povedať, že táto detekcia dáva zmysel.



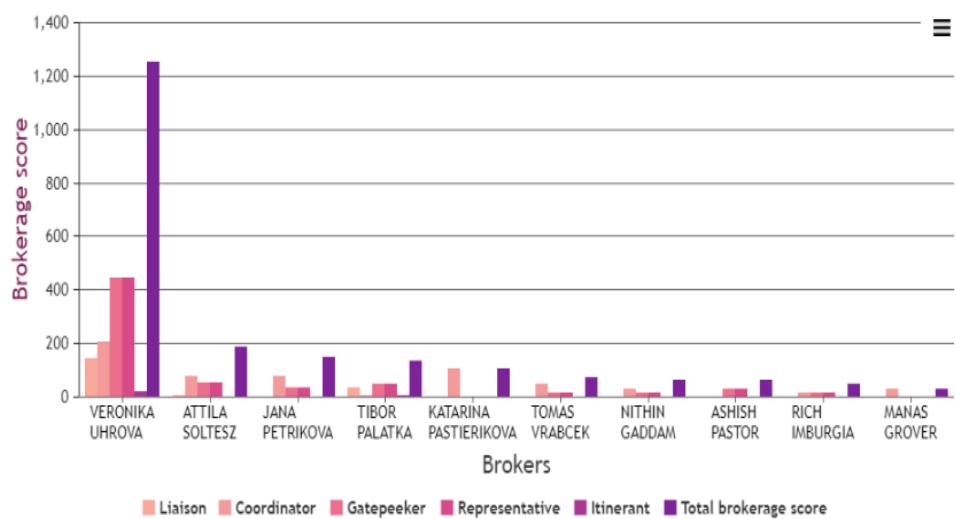
Obr. 43: Analýza jednotlivca - detail detekcie *leader* roly

### 8.2.5.2 Brokerage

V na obr. 44 a 45 je zobrazených desať najväčších *broker* aktérov v sieti spolu s grafom ich čiatkovým skóre pre každú *brokerage* rolu, ako aj celkové *brokerage* skóre.

TOP 10 BROKERS							
	Name	Coordinator	Itinerant	Gatekeeper	Representative	Liaison	Total
1	VERONIKA UHROVA	206	20	442	442	142	1252
2	ATTILA SOLTESZ	78	0	52	52	4	186
3	JANA PETRIKOVA	76	0	34	34	2	146
4	TIBOR PALATKA	6	4	46	46	32	134
5	KATARINA PASTERIKOVA	104	0	0	0	0	104
6	TOMAS VRABCEK	46	0	12	12	0	70
7	NITHIN GADDAM	30	0	16	16	0	62
8	ASHISH PASTOR	0	0	30	30	0	60
9	RICH IMBURGIA	14	0	16	16	0	46
10	MANAS GROVER	30	0	0	0	0	30

Obr. 44: Desať aktérov s najväčším *brokerage* skórom



Obr. 45: Desať aktérov s najväčším *brokerage* skórom - graf

## **9 Záver**

### **9.1 Možnosti rozšírenia a zdokonalenia práce**

## Literatura

- [1] J. Diesner, T. L. Frantz, and K. M. Carley, “Communication networks from the enron email corpus “it’s always about the people. enron is no different”,” *Computational & Mathematical Organization Theory*, vol. 11, no. 3, pp. 201–228, 2005.
- [2] X. Fu, S.-H. Hong, N. S. Nikolov, X. Shen, Y. Wu, and K. Xuk, “Visualization and analysis of email networks,” in *Visualization, 2007. APVIS’07. 2007 6th International Asia-Pacific Symposium on*, pp. 1–8, IEEE, 2007.
- [3] A. Chapanond, M. S. Krishnamoorthy, and B. Yener, “Graph theoretic and spectral analysis of enron email data,” *Computational & Mathematical Organization Theory*, vol. 11, no. 3, pp. 265–281, 2005.
- [4] G. Tang, J. Pei, and W.-S. Luk, “Email mining: tasks, common techniques, and tools,” *Knowledge and Information Systems*, vol. 41, no. 1, pp. 1–31, 2014.
- [5] A. Abnar, M. Takaffoli, R. Rabbany, and O. R. Zaïane, “Ssrm: structural social role mining for dynamic social networks,” *Social Network Analysis and Mining*, vol. 5, no. 1, p. 56, 2015.
- [6] S. Zehnalova, Z. Horak, and M. Kudelka, “Email conversation network analysis: Work groups and teams in organizations,” in *Advances in Social Networks Analysis and Mining (ASONAM), 2015 IEEE/ACM International Conference on*, pp. 1262–1268, IEEE, 2015.
- [7] “Do millennial and gen z consumers still use email?.” <https://www.bluecore.com/blog/do-millennials-use-email/>. 2017-03-30.
- [8] M. E. Newman, “The structure and function of complex networks,” *SIAM review*, vol. 45, no. 2, pp. 167–256, 2003.
- [9] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, “Fast unfolding of communities in large networks,” *Journal of statistical mechanics: theory and experiment*, vol. 2008, no. 10, p. P10008, 2008.
- [10] N. Crossley, E. Bellotti, G. Edwards, M. G. Everett, J. Koskinen, and M. Tranmer, *Social network analysis for ego-nets: Social network analysis for actor-centred networks*. Sage, 2015.
- [11] R. V. Gould and R. M. Fernandez, “Structures of mediation: A formal approach to brokerage in transaction networks,” *Sociological methodology*, pp. 89–126, 1989.
- [12] P. V. Marsden, “Brokerage behavior in restricted exchange networks,” *Social structure and network analysis*, vol. 7, no. 4, pp. 341–410, 1982.

- [13] R. S. Burt, *Structural holes: The social structure of competition*. Harvard university press, 2009.
- [14] K. Stovel and L. Shaw, “Brokerage,” *Annual Review of Sociology*, vol. 38, pp. 139–158, 2012.
- [15] E. S. Spiro, R. M. Acton, and C. T. Butts, “Extended structures of mediation: Re-examining brokerage in dynamic networks,” *Social Networks*, vol. 35, no. 1, pp. 130–143, 2013.
- [16] R. DeJordy and D. Halgin, “Introduction to ego network analysis,” *Boston MA: Boston College and the Winston Center for Leadership & Ethics*, 2008.
- [17] R. S. Burt, “Structural holes and good ideas,” *American journal of sociology*, vol. 110, no. 2, pp. 349–399, 2004.
- [18] S. P. Borgatti, “Structural holes: Unpacking burt’s redundancy measures,” *Connections*, vol. 20, no. 1, pp. 35–38, 1997.
- [19] “Active Record vs. Repository pattern repository pattern.” <https://www.rarous.net/weblog/271-active-record-vs-repository-pattern.aspx>. 2018-01-15.
- [20] “TeamNET data.” <http://inflex.cz:8075/TeamNETdata>. 2017-09-30.