

# VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

BRNO UNIVERSITY OF TECHNOLOGY

FAKULTA INFORMAČNÍCH TECHNOLOGIÍ  
ÚSTAV POČÍTAČOVÉ GRAFIKY A MULTIMÉDIÍ

FACULTY OF INFORMATION TECHNOLOGY  
DEPARTMENT OF COMPUTER GRAPHICS AND MULTIMEDIA

## DATAMINING Z JABBERU

BAKALÁŘSKÁ PRÁCE

BACHELOR'S THESIS

AUTOR PRÁCE

AUTHOR

JAROSLAV SENDLER

BRNO 2011



**VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ**  
BRNO UNIVERSITY OF TECHNOLOGY



**FAKULTA INFORMAČNÍCH TECHNOLOGIÍ**  
**ÚSTAV POČÍTAČOVÉ GRAFIKY A MULTIMÉDIÍ**

FACULTY OF INFORMATION TECHNOLOGY  
DEPARTMENT OF COMPUTER GRAPHICS AND MULTIMEDIA

## **DATAMINING Z JABBERU**

DATAMINING FROM JABBER

### **BAKALÁŘSKÁ PRÁCE**

BACHELOR'S THESIS

### **AUTOR PRÁCE**

AUTHOR

JAROSLAV SENDLER

### **VEDOUCÍ PRÁCE**

SUPERVISOR

Ing. JOZEF MLÍCH

BRNO 2011

**SEM VLOZTE  
ZADANI**

**PRACE**

## Abstrakt

Tato bakalářská práce se zabývá získáváním a analýzou dat ze sítě pro rychlé zasílání zpráv Jabber. Je zde popsán návrh klienta, který sbírá data a následně je ukládá do databáze. Tento klient plní funkci robota, který vhodně reaguje na zvolené dotazy. Dále jsou nashromážděné informace analyzovány. Pro tento účel byla nejprve provedena jejich transformace, která zohledňuje časové vlastnosti dat. Následně byla využita shluková analýza pomocí algoritmu  $k$ -means a nástroje RapidMiner. Výsledky získané procesem data mining jsou poté vhodně upraveny a prezentovány v grafické formě i pomocí tabulky.

## Abstract

This bachelor's thesis deals with the collection and analysis of data from the network for the Jabber instant messaging. There is a client described that collects data and then stores them in a database. The client performs the function of a robot that responds appropriately to the selected questions. Further information is collected for analysis. For this purpose their transformation was first made, that reflects the temporal characteristics of data. It was subsequently used algorithm  $k$ -means of tool RapidMiner. The results are presented by graphical form and using a spreadsheet.

## Klíčová slova

Jabber, XMPP, robot, data mining, shlukování,  $k$ -means, dolování z dat, transformace dat, RapidMiner.

## Keywords

Jabber, XMPP, robot, data mining, clustering,  $k$ -means, data transformation, RapidMiner.

## Citace

Jaroslav Sendler: Datamining z jabberu, bakalářská práce, Brno, FIT VUT v Brně, 2011

# Datamining z jabberu

## Prohlášení

Prohlašuji, že jsem tuto bakalářskou práci vypracoval samostatně pod vedením pana Ing. Jozefa Mlícha. Uvedl jsem všechny literární prameny a publikace, ze kterých jsem čerpal.

.....

Jaroslav Sendler

16. května 2011

## Poděkování

Tímto bych chtěl poděkovat mému vedoucímu bakalářské práce Ing. Jozefovi Mlíchovi za ochotu a kladný přístup při konzultacích. Dále za poskytnutí hardware na němž běžel program a sbíral data.

© Jaroslav Sendler, 2011.

*Tato práce vznikla jako školní dílo na Vysokém učení technickém v Brně, Fakultě informačních technologií. Práce je chráněna autorským zákonem a její užití bez udělení oprávnění autorem je nezákonné, s výjimkou zákonem definovaných případů.*

# Obsah

<b>1 Úvod</b>	<b>5</b>
<b>2 XMPP</b>	<b>7</b>
2.1 Architektura . . . . .	7
2.2 XML . . . . .	10
2.3 Stanza . . . . .	11
2.4 Rozšíření . . . . .	13
<b>3 Data mining</b>	<b>16</b>
3.1 Transformace dat . . . . .	18
3.2 Shlukování . . . . .	21
3.3 Programy . . . . .	25
<b>4 Implementace</b>	<b>27</b>
4.1 Architektura . . . . .	27
4.2 Návrh databáze . . . . .	28
4.3 Návrh robota . . . . .	30
4.4 Pokračování práce . . . . .	31
<b>5 Vyhodnocení výsledků</b>	<b>33</b>
5.1 Manuální rozbor dat . . . . .	33
5.2 Programové vyhodnocení . . . . .	34
<b>6 Závěr</b>	<b>38</b>
<b>A Obsah CD</b>	<b>42</b>
<b>B Slovník zkratk</b>	<b>43</b>
<b>C Stanza - základní schéma</b>	<b>44</b>
<b>D Návrh databáze</b>	<b>48</b>
<b>E Přehled příkazů robota</b>	<b>49</b>
<b>F Přehled vztahů uživatelů</b>	<b>50</b>

# Seznam obrázků

2.1	Distribuovaná architektura Jabber. . . . .	8
2.2	Rozebraná struktura Jabber ID. . . . .	9
3.1	Proces dobývání znalostí z databází podle knihy autora Fayyad [4]. . . . .	17
3.2	Srovnání výpočtu vzdáleností od bodu $x_1$ [2]. . . . .	22
3.3	Algoritmus $k$ -means zobrazený pomocí konečného automatu. . . . .	24
4.1	Struktura architektury bakalářské práce. . . . .	28
4.2	Vybraná část struktury databáze. . . . .	29
4.3	Vybraná část struktury robota. . . . .	31
5.1	Struktura procesu v nástroji RapidMiner. . . . .	35
5.2	Výsledné shluky převedeny do 3D prostoru. . . . .	35
5.3	Průběh statusů shluků podle času. . . . .	36
C.1	Ukázka „šíření“ <i>User Tune</i> . . . . .	45
D.1	Struktura databáze . . . . .	48

# Seznam tabulek

2.1	Přehled Jabber serverů. . . . .	10
3.1	Ukázka tabulky <i>presence</i> . . . . .	20
3.2	Ukázka modifikované tabulky <i>presence_modify</i> . . . . .	20
5.1	Přehled manuálního rozboru dat. . . . .	34
5.2	Procentuální shoda jednotlivých uživatelů. . . . .	37
E.1	Přehled příkazů pro robota. . . . .	49
F.1	Zobrazení uživatelského JID na čísla. . . . .	50
F.2	Procentuální shoda jednotlivých uživatelů. . . . .	51



# Příklady

2.1	Ukázka základního XML dokumentu. . . . .	11
2.2	Použití elementu <i>message</i> . . . . .	11
2.3	Použití elementu <i>iq</i> . . . . .	12
2.4	Použití elementu <i>presence</i> . . . . .	12
2.5	Začátku vysílání rozšířeného statusu. . . . .	13
2.6	Dotaz na podporované protokoly. . . . .	14
3.1	Metoda <i>k</i> -means byla převzata z [2]. . . . .	25
C.1	Popis elementu <i>iq</i> . . . . .	44
C.2	Popis elementu <i>message</i> . . . . .	44
C.3	Popis elementu <i>presence</i> . . . . .	45
C.4	Informování serveru o právě přehrávající hudbě. . . . .	46
C.5	Server informuje uživatele podporující rozšíření o stavu <i>user@jabber.com</i> . . . . .	46
C.6	Server přepośle informace o přehrávané hudbě všem otevřeným spojením uživatele <i>user@jabber.com</i> . . . . .	47
C.7	Uživatel ukončil „vysílání“ rozšířených zpráv o svém stavu. . . . .	47
C.8	Server informuje klienty o ukončení šíření rozšířeného statusu. . . . .	47

# Kapitola 1

## Úvod

V současné době jsou řazeny k nejrozšířenějším psaným dorozumívacím prostředkům komunikační sítě pracující v reálném čase. Samozřejmě je tento jev k vidění až koncem minulého desetiletí. Psaná komunikace je již několik století využívána jako prostředek k dorozumívání lidí mezi sebou. Za toto období prošla výrazným pokrokovým vývojem, proto lze u ní nalézt mnoho časových mezníků, které ji výrazně ovlivnily. Jako příklad jednoho z prvních komunikačních kanálů lze uvést posly, kteří často nesly zprávy na vzdálenosti několika kilometrů. Dalším výrazným prvkem ve vývoji dorozumívacích prostředků bylo zavedení pošty a objevení telegrafu.

Příchodem internetu nastal v komunikaci zásadní zlom. Postupným rozšířením pokrytí a dostupnosti této technologie začalo vznikat mnoho nových komunikačních prostředků. Příkladem mohou být elektronické zprávy, RSS zprávy nebo komunikační sítě, které pracují v reálném čase. Právě poslední zmíněná metoda zažívala v moderní době velký růst v oblasti popularity, ať už v podobě ICQ protokolu nebo dnes velmi rozšířené sociální sítě Facebook. Jedna z hlavních výhod těchto služeb, například v porovnání s klasickou poštou, je jejich rychlost doručení. Naproti běžné elektronické poště jsou uživatelé informováni o stavu příjemce zprávy. Real-time komunikační prostředky nabízejí služby, kterými je možné zjistit, zda se příjemce zprávy nachází u dorozumívacího zařízení. Díky této schopnosti je uživateli umožněno zasílat zprávy a obratem na ně očekávat odpovědi.

S přibývajícími elektronickými daty, na poli ať už vědních nebo praktických oborů, přichází potřeba tyto informace uchovávat. K těmto účelům slouží databáze, které se svou strukturou zaměřují především na snadnou kontrolu dat, vyhledávání a jejich analyzování. S rostoucím obsahem se ale uložené informace stávají pouze daty bez významu. Proto vznikla nová disciplína nazvaná *data mining*, která si klade za cíl znovunalezení „ztracených“ informací (na první pohled neviditelných) nebo nalezení informací úplně nových.

Předmětem této bakalářské práce je seznámení se s problematikou komunikace probíhající přes Jabber síť. Konkrétním cílem je vytvoření jednoduchého Jabber klienta, který by byl schopen získávat statistická data. Nashromážděná data jsou dále využita pro pozdější analýzu a grafickou reprezentaci informací z nich získaných. Hlavním cílem této práce je získat neznámé informace z komunikační sítě Jabber.

Technická zpráva je tvořena z šesti kapitol. Kapitola **druhá** je tvořena popisem protokolu XMPP, na kterém je postavena komunikační služba Jabber pracující v reálném čase. Je zde popsána architektura sítě a základní stavební kameny XMPP protokolu, které jsou využívány sítí Jabber jako nástroj k zprostředkování jednotlivých služeb. V závěru tohoto oddílu je věnována část vybraným standardům, které rozšiřují základní XMPP protokoly. Jsou zde uvedena pouze ta rozšíření, která byla po konzultaci s vedoucím práce, označena

za relevantní k této práci.

Obsah **třetí** kapitoly je zaměřen na popis procesu data mining. Zabývá se začleněním této metody do komplexnějšího procesu, který je nazýván získávání znalostí z databází. Další součástí této kapitoly se zabývá nezbytnými kroky transformace dat, která jsou získána z Jabber komunikace. Při tomto procesu jsou data převedena do vhodné podoby pro proces dolování z dat. Následuje část, kterou jsou popsány vybrané metody dolování dat a také je zde podrobně rozebrán algoritmus  $k$ -means. Tento algoritmus je využit pro samotné dolování z dat, které je zprostředkováno aplikací RapidMiner. Popis tohoto programu a dalších vybraných nástrojů uzavírá třetí kapitolu.

Implementační část této práce je popsána kapitolou **čtvrtou**. Je členěna na oddíly, které odpovídají vybraným částem architektury této práce. Elementy architektury, které jsou popsány vlastní podkapitolou, jsou rozšířeny o zjednodušený návrh struktury v podobě grafických schémat. Závěr této kapitoly se zamýšlí nad rozšířeními robota a nad dalším pokračováním této práce.

Výsledky získané touto prací jsou prezentovány v kapitole **páté**. Je zde ukázán model používaný programem RapidMiner pro dobývání znalostí. Konkrétně jde o shlukování uživatelů sítě Jabber do skupin podle času stráveném v samotné síti a jejich stavu. Výsledky jsou představeny pomocí různých grafických entit, ať už v podobě diagramu nebo grafu, který je transformovaný do dvourozměrného prostoru.

Nemalý informativní obsah tvoří i **přílohy**. V první řadě je to slovník zkratk, shrnující slovní spojení z celého tohoto dokumentu. Následuje podrobnější popis základních entit stanzy a přehled průběhu rozšíření. V další části je ukázán kompletní návrh databáze, který doplňuje základní popis ze čtvrté kapitoly. Přehled příkazů, na které robot vytvořený v této práci reaguje, je ukázán v poslední části příloh.

## Kapitola 2

# XMPP

V následující kapitole jsou, pro usnadnění a jednodušší pochopení, rozebrány základní stavební kameny protokolu Extensible Messaging and Presence Protocol (XMPP). Konkrétně jsou zde popsány stávající vlastnosti implementace, architektura protokolu XMPP obecně [23, 24] a další detaily protokolu [1, 25, 14]. Vzhledem k požadavkům na dolování v datech popsaných v následující kapitole je kladen důraz na vybraná rozšíření [22, 12]. Tato rozšíření tvoří základ pro některé rozšířené statusy, jako je například User Tune [18], User Mood [21], User Location [6] a další. Další informace použité pro popis a pochopení XML jazyka byly čerpány z [10, 9].

Vznik samotného protokolu XMPP je datován do roku 2004 (březen), kdy na něj byl přejmenován Jabber. Původní projekt Jabber byl vytvořen roku 1998 autorem Jeremie Millerem, který ho založil za účelem vytvořit svobodnou otevřenou IM službu. Uvedený projekt měl obsahovat tři základní vlastnosti, do kterých se zahrnují jednoduchost a srozumitelnost pro implementaci, jednoduchost v oblasti šíření a otevřenost podobě veřejně dostupného popisu samotného protokolu. Základní vlastnosti a výhody klientů a serverů budou podrobněji popsány níže. Roku 1999, 4. ledna byl vytvořen první server se jménem Jabber. Komunita vývojářů se chopila iniciativy a vytvořila klienty, kteří dokázali se serverem komunikovat, pro různé platformy (Linux, Macintosh, Windows). Roku 2004 byl protokol XMPP přidán mezi RFC<sup>1</sup> dokumenty. Základní norma popisující obecnou strukturu protokolu je RFC 3920 [23] a RFC 3921 [24], který se zaměřuje na samotný instant messaging a zobrazení stavu. Další zdokumentovaná rozšíření jsou vydávána v podobě tzv. XEP (XMPP Extension Protocol) dokumentů, které jsou známé také pod starším názvem JEP (Jabber Enhancement Proposal). Dnešní počet těchto norem se blíží k číslu 300. Každý XEP obsahuje stav vývoje (schválení), ve kterém se zrovna nachází.

Jako bezpečnostní prvky jsou zde podporovány SASL, TLS a GPG. XMPP protokol je postaven na obecném značkovacím jazyce XML, proto vlastnosti popsané dále v této kapitole platí i pro tento protokol.

### 2.1 Architektura

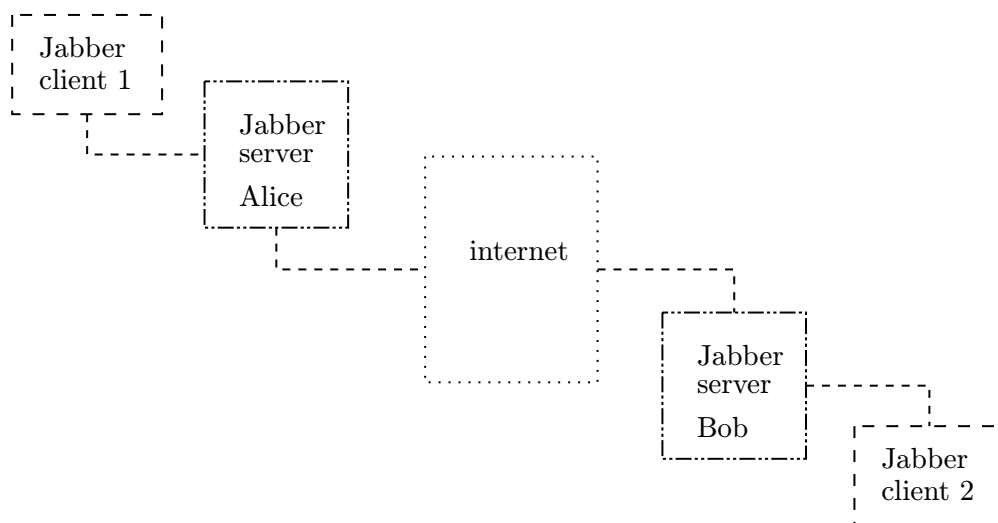
Dobře navržená internetová technologie je tvořena správně fungujícími komponenty, které mezi sebou dokáží vytvořit spojení a následně započít komunikaci. Pro popis Jabber architektury v této práci bylo čerpáno z [1, 25]. Tato struktura se nejvíce podobá struktuře posílání e-mailů. Hlavní předností Jabber sítě je, tak jako u elektronické pošty, její decentra-

---

<sup>1</sup>RFC request of comments — žádost o komentáře

lizace. V případě Jabberu je decentralizace chápána jako možnost provozovat vlastní server, na rozdíl od jiných komunikačních systémů jako je například Facebook, kde existuje pouze jediný poskytovatel služby. V případě serveru je kladen důraz na spolehlivost a rozšiřitelnost a u klienta na uživatele. Každý server pracuje samostatně, což znamená, že chod ani výpadek jiné datové stanice žádným způsobem jeho běh neovlivní. V případě výpadku jiného serveru bude nedostupný pouze seznam kontaktů a služeb, které registrovaným uživatelům poskytoval.

Obrázek 2.1 znázorňující distribuovanou architekturu Jabberu byl převzat z [1] a doplněn o názvy jednotlivých komponent. Komunikace dvou Jabber klientů probíhá za účasti jejich serverů a sítě, která je spojuje. Spojení mezi nimi bývá často šifrováno.



Obrázek 2.1: Distribuovaná architektura Jabber.

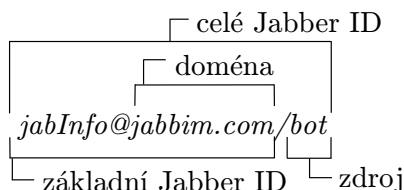
Architektura Jabber serverů využívá velké množství mezi-doménových připojení podobně jako internetový systém elektronické pošty. Komunikace klienta z jedné domény s klientem z jiné na rozdíl od e-mailového modelu nevyžaduje spolupráci třetích stran. Klient se spojí s „domácím“ serverem, který přímo naváže spojení se serverem požadovaného klienta. Tyto vlastnosti jsou zárukou pro bezpečný přenos zpráv, znemožňující „krádeže“ JID<sup>2</sup>, který je popsán níže, a spamování.

## Jabber ID

Jabber ID (JID) je jednoznačný virtuální identifikátor uživatele na síti. V případě založení účtu nejsou rozlišována velká a malá písmena, což znamená, že Jabber není case-sensitive. Jednoznačný Jabber identifikátor je složen ze dvou částí: *Jabber bare* neboli čisté ID a *resource* [23]. Základní část na první pohled připomíná e-mailovou adresu *user@server*. Druhá část slouží k přesné identifikaci jednotlivých spojení. Je použita ke směrování síťového provozu s uživateli v případě otevření většího množství spojení pod jedním uživatelem. Společně Jabber bare a resource tvoří tzv. *full JID* — *user@server/resource* například *jabInfo@jabim.cz/bot*. Jednotlivé části uživatelského jména popsané v tomto odstavci jsou ukázány v obrázku 2.2.

<sup>2</sup>uživatelské jméno

Další vymoženost JabberID oproti e-mailové adrese je jeho možnost používat prakticky libovolné národní znaky u doménových jmen a uživatelských účtů [25]. Využíváním kó-



Obrázek 2.2: Rozebraná struktura Jabber ID.

dování UNICODE, se XMPP stává plně mezinárodní a není jako jiné protokoly omezen rozsahem ASCII tabulky. Přestože je tato vymoženost k dispozici, doposud není žádným výrazným způsobem využívána.

## Klient

Klient je často jednoduchá aplikace pracující se vzdálenými službami, které jsou provozovány serverem. V této práci je zastoupen robotem s konzolovým rozhraním. XMPP svou architekturou nutí, aby byl co nejjednodušší. Vlastnosti, které by měl mít, jsou shrnuty podle [14] do tří bodů:

1. komunikace s jedním Jabber serverem pomocí TCP socketu, který garantuje spolehlivé doručení zpráv na rozdíl od UDP. Nad tímto transportním protokolem dále běží kryptografický protokol TLS, který zabezpečuje komunikaci klient-server a server-server.
2. rozparsování a následná interpretace příchozí XML zprávy „stanza“ (kapitola 2.3).
3. porozumění sadě zpráv (*message*, *iq*, *presence*) z Jabber jádra [23].

## Server

Informace použité pro popis XMPP serveru byly čerpány z [14]. K hlavním charakteristikám serveru oproti klientovi, jehož základní vlastností byla jednoduchost, patří stabilita a bezpečnost. Je pro něj vyhrazen TCP port 5222. Komunikace mezi servery je realizována přes port 5269. Každý server uchovává seznam zaregistrovaných uživatelů, který nevykazuje žádný jiný server. Zaregistrovaní uživatelé v daném seznamu se mohou do sítě připojovat pouze přes něj. To zajišťuje nemožnost „krádeže“ účtu. Protože XMPP komunikace probíhá přes síť, musí mít každá entita adresu, v tomto případě nazvána JabberID. XMPP spoléhá na DNS což znamená, že používá jména na rozdíl od IP protokolu.

Server Jabber je systém spravující tok dat mezi jednotlivými komponentami, které společně tvoří Jabber služby. Například *Jabber Session Manager* (JSM) poskytne funkce pro IM komunikaci a práci se seznamem kontaktů. Komunikace mezi jednotlivými servery, jak je uvedeno na obrázku 2.1, je zprostředkována za pomoci komponenty *S2S* (server to server). Při připojení klienta k serveru je komunikace řízená pomocí *C2S* (client to server). Jak již bylo řečeno, Jabber síť využívá doménová jména místo špatně zapamatovatelných IP adres. Pro tento způsob identifikace je určena služba *dnsrv*, která se stará o překlad názvů. V podstatě je to komponenta, která zajišťuje směrování paketů na jiný server.

V tabulce 2.1 jsou shrnuty informace o serverech Jabberu. První sloupec tvoří jméno, následuje programovací jazyk, v němž je napsán. Většina aplikací pro servery je vydávána pod licencí GPL<sup>3</sup>. U všech aplikací byla zkoumána nejaktuálnější verze. Její číslo lze nalézt ve třetím sloupci. Všechny servery lze provozovat na operačním systému Linux a Windows. Na platformě Mac OS mohou být použity všechny zde jmenované vyjma jabberd2. Pět z šesti zde představených programů pro server Jabber jsou stále vyvíjeny, tedy kromě jabberd14. Hlavním účelem tabulky je prezentovat důležité vlastnosti serverů v oblasti podpory rozšířených statusů. Jedná se o standardy *pubsub*<sup>4</sup> (XEP-0060) [12] a o jeho verzi, která je více zaměřena na uživatele *pep*<sup>5</sup> (XEP-0163) [22]. Obě tato rozšíření tvoří nezbytnou základnu pro *rozšířené statusy* a proto je jejich podpora jak u serverů, tak klientů vyžadována. Podrobněji toto téma bude rozebráno v některé následující podkapitole.

Server	Jazyk	Verze	XEP-0060	XEP-0163
ejabberd	Erlang/ Top	2.1.6	ANO	ANO
Openfire	java	3.6.4	ANO	ANO
jabbred2	c	2.2.11	NE	NE
jabbred14	c, c++	1.6.1.1	ANO	NE
Prosody	lua	0.7.0	NE	ANO
Tigase	java	5.0.0	ANO	ANO

Tabulka 2.1: Přehled Jabber serverů.

Z výše uvedené tabulky je zřejmé, že aplikace pro servery, které jsou stále ve vývoji, podporují tzv. *rozšířené statusy*. Tedy kromě programu jabbred2.

## 2.2 XML

Jazyk XML (eXtensible Markup Language) [9], metajazyk pro deklaraci strukturovaných dat, je jádrem protokolu XMPP. Samotný jazyk vznikl rozšířením metajazyka SGML, jež slouží pro deklaraci různých typů dokumentů. Základní vlastností je jednoduchá definice vlastních značek (tagů). Dokument XML se skládá z elementů, které můžeme navzájem zanořovat. Vyznačujeme je pomocí značek — počáteční a ukončovací. Pomocí tohoto jazyka je tvořena *stanza* popsána v následující kapitole.

Ukázka možné struktury dokumentu psaného jazykem XML je zobrazena na příkladu 2.1. Standardně je předpokládáno, že je psán v kódování UTF-8 [10], ale je-li jako v tomto případě použito jiné, musí být konkrétní kódování uvedeno na jeho počátku. V opačném případě nemusí být obsah správně zobrazen. Na začátku dokumentu se také uvádí verze XML, ve které je dokument psán (1. řádek příkladu). Následuje kořenový element, který je uzavřen na samotném konci dokumentu. 4. řádek prezentuje možnost použití prázdného elementu, který obsahuje jeden atribut s názvem zkratky fakulty. Velký význam zde mají úhlové závorky. Jsou jimi z obou stran obaleny všechny elementy.

<sup>3</sup>General Public License — všeobecná veřejná licence GNU

<sup>4</sup>Publish-Subscribe

<sup>5</sup>Personal Eventing Protocol

---

```

1      <?xml version="1.0" encoding="iso-8859-2"?>
2      <fakulta>
3          <název>Fakulta informačních technologií</název>
4          <zkratka fakulty="FIT"/>
5          <typy studia>
6              <bakalářské titul="Bc."></bakalářské>
7              <magisterské></magisterské>
8              <doktorské></doktorské>
9          </typy studia>
10     </fakulta>

```

---

Příklad 2.1: Ukázka základního XML dokumentu.

## 2.3 Stanza

Základní jednotkou pro komunikaci založenou na XML je stanza. Z jednoduššího pohledu je možné se na ni dívat jako na jeden dlouhý XML soubor. Při zahájení komunikace se tento soubor „otevře“. Jeho samotné uzavření probíhá až při odhlášení od sítě, neboli přepnutí klienta do stavu offline. Stanzu je tedy možné vnímat jako stream, který obsahuje všechna data probíhající komunikace. Mezi elementy používané pro komunikaci klienta se serverem patří tyto tři: *message*, *presence* a *iq*. Každý zde uvedený člen má svůj jednoznačný význam. V následujících odstavcích jsou jednotlivé části stanzy blíže definovány a na reprezentativních příkladech jsou ukázány jejich základní struktury a možnosti využití v praxi.

První prvek, který bude charakterizován je označen anglickým výrazem *message* (zpráva). Jak již název napovídá, slouží k posílání zpráv všeho druhu. Je to základní metoda pro rychlý přenos informací z místa na místo. Zprávy jsou typu „push“, což znamená, že jsou odeslány a není očekávána žádná aktivita od příjemce, která by přijetí potvrdila. Jedno z dosavadních využití se nachází v klasické komunikace po internetu, tzv. instant messaging (IM). K dalším možným použitím patří skupinový chat a oznamovací nebo upozorňující zprávy. Každá z těchto zpráv je tvořena z minimální povinné struktury. Tak jako u klasické poštovní korespondence nesmí chybět adresa odesílatele a adresa příjemce, kterému je zpráva adresována. Podle možnosti použití jsou zprávy děleny do kategorií. Jmenovitě toto rozdělení implementuje atribut *type*, který může nabývat jednu ze čtyř hodnot. Jsou rozlišovány zprávy pro komunikaci mezi dvěma entitami, skupinový chat, upozornění, chybová zpráva a v neposlední řadě zpráva bez kontextu vyžadující odpověď příjemce. Nakonec nesmí být opomenut blok zprávy, pro uživatele IM nejdůležitější, nesoucí vlastní obsah.

Základní použití struktury elementu *message* je prezentováno na příkladu 2.2. Na prvním řádku je uveden atribut, značící odesílatele. Druhý řádek obsahuje JID klienta, který zprávy přijímá. Následuje informace o typu zprávy a poté je uveden element *body* nesoucí samotný obsah.

---

```

1      <message from="user@jabber.com"
2              to="jabinfo@jabber.com/bot"
3              type="chat"
4          <body> Kolik je hodin? </body>
5      </message>

```

---

Příklad 2.2: Použití elementu *message*.

Další částí stanzy je poskytována struktura pro *request-response* (žádost-odpověď) vazbu, podobnou metodám GET, POST a PUT z protokolu HTTP [25]. Zkráceně je označována



pomocí dvou počátečních písmen *Info/Query* neboli IQ. Na rozdíl od elementu *message* tvoří *iq* spolehlivější přenos, optimalizovaný pro výměnu dat (binární data). K dalším rozdílům patří povinnost příjemce odpovědět na každou přijatou zprávu, neboli potvrdit její doručení. Skutečnost, že je na právě požadovanou zprávu odpovězeno, zajišťuje parametr *id*. Iq dotaz nebo odpověď musí obsahovat stejnou hodnotu tohoto atributu jako zpráva vytvořená žádajícím subjektem. Další povinný atribut rozděluje iq na čtyři typy. Jednotlivé žádosti na proces nebo akci jsou posílány samostatně [24]. V příloze C je uvedena rozsáhlejší struktura tohoto elementu. Použití nachází v případech, které nastavují, žádají nebo informace posílají. Tato struktura je využívána pro novou registraci, posílání seznamu kontaktů a další.

Příklad 2.3 znázorňuje základní použití elementu *iq*. Uživatel *user* posílá dotaz na získání seznamu kontaktu (řádek 5.).

---

```

1      <iq from="user@jabber.com/doma"
2          to="user@jabber.com"
3          id="uhhfw23648"
4          type="get"
5      <query xmlns="jabber:iq:roster" />
6      </iq>

```

---

Příklad 2.3: Použití elementu *iq*.

Poslední a pro tuto práci nejdůležitější prvek stanzy je *presence*. V případě, že nemá určeného příjemce, tak funguje způsobem jako broadcast. Což znamená, že jsou informace směřovány všem klientům, kteří jsou zaregistrováni k jejímu odběru. Presence v českém překladu informace o stavu (přítomnost) rozesílá dostupnost ostatních entit v síti. Jedná se tedy o nastavení uživatelské dostupnosti tak jako na jiných real-time komunikačních a sociálních systémech.

Existuje několik základních stavů statusů, které reprezentují aktuální dosažitelnost uživatele. Tento jev je vyjádřen pomocí elementu *show*, který disponuje čtyřmi možnostmi. První oznamuje, že je uživatel k dispozici a schopen aktivní komunikace. Druhá často se vyskytující možnost naznačuje, že je subjekt krátkou dobu pryč od svého IM klienta. Tento a další dva stavy, popsané dále, jsou často změněny bez lidského zásahu (pomocí PC nebo jiného zařízení) prostřednictvím funkce známé jako „auto-away“. Poslední dva stavy charakterizují delší časové období nečinnosti. Tato oznámení o změně stavu uživatele jsou často zasílána pouze kontaktům, které se nacházejí v režimu online. Tato optimalizace přispívá ke snížení síťového provozu, jelikož presence v reálném čase při komunikaci využívá velké množství šířky pásma.

Základní použití *presence* je zobrazeno v příkladu 2.4. Kontakt *jabinfo@jabber.com/bot* (1. řádek) posílá informace o svém stavu (řádek č. 2) a svůj status (č. 3).

---

```

1      <presence from="jabinfo@jabber.com/bot"
2          <show> online </show>
3          <status> Jsme zde. </status>
4      </presence>

```

---

Příklad 2.4: Použití elementu *presence*.

Obsáhlejší struktura elementu *presence* je zobrazena v příloze C, kde je rovněž k nalezení přehled všech možných stavů.

Jak již bylo zmíněno v části o Jabber ID, Jabber podporuje práci s více současně připojenými klienty k jednomu Jabber účtu. Vysvětlení funkčnosti bude prezentováno na příkladu uživatele přihlášeného na stolním počítači a z klienta v mobilním telefonu. U obou těchto připojení je použit stejný Jabber bare, ale odlišného resource, například *domov* a *mobile*. Právě tento rozdíl v tzv. „full“ adrese účtu zajišťuje jednu ze dvou možných podmínek pro správnou adresaci zpráv. Druhá možnost, která bude uplatněna při použití adresy účty pouze ve formě Jabber bare, je nastavení priority u jednotlivých programů. Priorita je číslo v rozsahu hodnot od -128 do 127, kde klient s větší prioritou má přednost před klientem s nižší. Nastane-li případ připojení více klientů se stejnou prioritou, každý server se při rozepisování zpráv zachová podle vlastní implementace. Některé rozešlou zprávy všem klientům, jiné naopak jen poslednímu přihlášenému.

## 2.4 Rozšíření

Dále se tato technická zpráva zabývá rozšířeními protokolu XMPP o další vlastnosti, k jejichž popisu slouží XEP. Pro tuto práci jsou nepostradatelné „statusy“, pro které tvoří základ standardy XEP-0060 [12] a XEP-0163 [22] zkráceně PEP<sup>6</sup>. Obě tato rozšíření umožňují strukturovaně pracovat, používat a přenášet další XEP protokoly. Jako příklady relevantní k práci jsou zde uvedeny protokoly *User Location* (kde se uživatel právě nachází) [6], *User Tune* (co uživatel poslouchá za hudbu) [18], *User Mood* (aktuální nálada uživatele) [21] a *User Activity* (co uživatel právě dělá) [11]. Jsou to tedy protokoly založené na PEP, které vyžadují podporu nejen v klientech, ale i na straně serveru (zobrazuje tabulka 2.1). S touto informací úzce souvisí další protokol XEP-0115 [7], který umožňuje zjistit podporované schopnosti klienta, případně, které informace je ochoten přijímat. Tato vlastnost bude popsána níže v části zabývající se podporovanými vlastnostmi.

Všechna tato rozšíření by mohla být přidána přímo do statusu viz příklad 2.4, avšak ten je primárně určen k informování o přítomnosti na IM síti. Hlavní rozdíl mezi PEP a obyčejným posílání stavu pomocí presence je v pravomoci klienta přijmout nebo odmítnout informaci, na rozdíl od presence, jež je přijata vždy.

Základ přenosu informací začíná na straně klienta, který chce všechny ve svém roster listu (seznam kontaktů), informovat o statusu. Zašle zprávu obalenou v elementu *iq* serveru. Ukázka této zprávy je prezentována na příkladu 2.5, který znázorňuje zaslání informace o druhu hudby, kterou v danou chvíli uživatel poslouchá. Využívá k tomu rozšíření *User Tune*, definovaném na řádku číslo 5. Základ zprávy oznamující začátek vysílání informací o rozšířených statusech je vždy stejný. Liší se pouze řádkem 3. a obsahem elementu *item* v příkladu 2.5.

---

```

1      <iq from='user@jabbim.com' type='set' id='pub1'>
2          <pubsub xmlns='http://jabber.org/protocol/pubsub'>
3              <publish node='http://jabber.org/protocol/tune'>
4                  <item>
5                      <tune xmlns='http://jabber.org/protocol/tune'>
6                          <artist>Daniel Landa</artist>
7                          <length>255</length>
8                      ...

```

---

Příklad 2.5: Začátku vysílání rozšířeného statusu.

---

<sup>6</sup>Personal Eventing via Pubsub

V případě úspěšného přijetí *iq* zprávy serverem, každý, kdo se zaregistroval k odebrání rozšířených statusů, obdrží oznámení ve formě *message*. Oznámení bude také doručeno všem *resources*. Celá zpráva i všechny další náležitosti jsou uvedeny v příloze C.

## Podporované vlastnosti

Jednotlivá rozšíření protokolu XMPP jsou nepovinná, a proto nemusí být ve všech klient-ských aplikacích podporována. Pro zjištění podporovaných rozšíření se používá XEP-0115 Entity Capabilities [7]. Toto rozšíření výrazně snižuje počet a velikost komunikací a přenosů zpráv mezi uživateli. Dotazem zobrazeným na příkladu 2.6 je zjištěna schopnost jednotlivých klientů, kterou následně server využije pro správné směrování rozšířených statusů. Všechny zde zmiňované rozšíření a protokoly z této kapitoly je možné u každého klienta (seznam klientů obsahuje tabulka v příloze C) vyčíst z atributu *ver* (druhá část u atributu *node*), který je vypočítán ze všech podporovaných protokolů klienta, viz [7].

---

```
1<iq from="user@jabbim.com" id="disco1"
2  to="jabinfo@jabbim.com/bot" type="get">
3  <query xmlns="http://jabber.org/protocol/disco#info"
4    node="http://code.google.com/p/exodus#QgayPKawpkPSDYmwT/WM94uAlu0=" />
5</iq>
```

---

Příklad 2.6: Dotaz na podporované protokoly.

## Další rozšíření

V následujících několika odstavcích budou přiblíženy specifikace jednotlivých rozšíření XEP, které slouží jako zdrojová data pro dolování a jsou relevantní k tématu práce.

Prvním rozšířením, nad rámec základních vlastností Jabberu, které zde bude podrobněji rozebráno, je elektronická verze klasické vizitky neboli *VCard*. Jeho specifikací se zabývají dva standardy. Jelikož novější verze XEP dokumentu [13] se v době psaní této práce nacházela ve stavu „experimental“, což znamená, že ještě není schválena jako standard, je pouze ve stavu návrhu. Proto bylo použito verze starší [16]. Jednoduše řečeno je *VCard* struktura, která nese informace o uživateli jako je jméno, příjmení, e-mail, adresa bydliště i zaměstnání a další údaje. Data jsou dále zveřejňována na síti, z čehož vyplývá, že jsou dostupná ostatním uživatelům. Vyplnění těchto osobních údajů je dobrovolné a tak se u některých uživatelů nachází pouze přezdívka a JID, které jsou často předdefinovány automaticky. Nedílnou součástí všech sociálních a komunikačních systémů jsou malé fotografie, loga nebo ikony, kterými se uživatelé prezentují. V síti Jabber tomu není jinak, a proto je samotný obrázek zahrnut přímo do *VCard* v položce *photo*. Podrobnější informace o jeho nastavení a přijímání je možné nalézt v *vCard-Based Avatars* [19], který jej definuje.

Díky základní podmínce XMPP protokolu (otevřenost) existuje mnoho různých aplikací, pomocí kterých lze v síti Jabber komunikovat. S programy, používanými uživateli, úzce souvisí další zde implementované rozšíření. Jedná se o realizaci *Software Version* dokumentu [17], který se právě zabývá získáváním informací o samotných aplikacích. Je-li toto rozšíření podporováno je díky němu možné zjistit jméno a verzi používané aplikace. Informace o operačním systému často nejsou kvůli bezpečnosti ani vyplněny. Podrobnější informace o softwarové výbavě klienta je možné zjistit pomocí XEP [7], o kterém již bylo dříve psáno v odstavci zabývajícím se podporovanými vlastnostmi klientských aplikací.

S rozšířením tzv. „chytrých“ mobilních zařízení mezi širší veřejnost vzniklo několik nových disciplín spojených s určováním zeměpisné polohy, jako je například geocaching. Geografická poloha je přenášena ve formě souřadnic popisující přímo zeměpisnou šířku a délku. Současně lze informaci o poloze přenášet i slovně ve formě adresy. Příkladem slovního popisu je ulice, číslo popisné, město a další. Mnoho aplikací, které mají k dispozici GPS přijímač, vysílají a aktualizují zeměpisné informace automaticky, například po určité době nebo změně polohy o určitou vzdálenost. Toto a další níže popsání rozšíření jsou postaveny na již zmiňovaném PEP. Některé části protokolů jsou zjednodušeny a připraveny tím pro „mobilní instant messaging“.

Pro sdělení informací o stavu klienta není v základní verzi Jabberu mnoho. Pomocí presence je možné „pouze“ prozradit, zda je uživatel připraven komunikovat nebo je momentálně nedostupný a to v několika verzích lišících se délkou nepřítomnosti. Pokročilejší nastavení statusu nabízí *User Mood* [21] a to ve formě sdělení současné nálady, jako je například radost. Další možné upřesnění činnosti uživatele jsou definovány v *User Activity* [11], kde každá činnost je složena z povinné obecné kategorie a nepovinné, která informaci upřesňuje. Příkladem může být *eating* a *having\_a\_snack* tj. uživatel jí, uživatel svačí.

K poslednímu rozšíření implementovanému v této práci patří *User Tune* [18], které umožňuje uživateli šířit informace o aktuálně poslouchané hudbě. Některé dnešní hudební přehrávače dokáží automaticky spolupracovat s IM klientem a předávat informace o hudbě bez nutného lidského zásahu. Ve zprávě jsou tedy přenášeny informace o skladbě, interpretovi, albu a další informace, které mohou být získávány z MP3 ID3v1 nebo novější ID3v2 tag.

Podpora výše popsaných rozšíření v aplikacích je poměrně malá. Například v předcházející zmiňované části o poslouchané hudbě, při stavu, kdy program toto rozšíření nepodporuje, je posíláno pomocí normální presence. Jméno skladatele, alba a další podrobnosti jsou shrnuty do statusu, tudíž jsou doručeny všem uživatelům ze seznamu kontaktů.

## Kapitola 3

# Data mining

Třetí kapitola se zabývá procesem dobývání znalostí z databází. Popisuje jej jako disciplínu, která vznikla za účelem vytěžení informací z dat, která jsou v nepřehledném množství ukládána v databázích. Díky velikosti dnešních disků objem ukládaných dat neustále roste. S tím také úzce souvisí zvětšující se poměr nepotřebných a zašumělých dat vůči užitečným informacím. V této kapitole jsou mimo jiné popsány metody používané k dolování z dat, které jsou relevantní k této práci.

Na začátku kapitoly je rozebrán pojem získávání znalostí databází, jehož jednu podstatnou část tvoří samotný data mining. Dále je vysvětlena základní terminologie, pro kterou bylo čerpáno z [8]. Cílem první podkapitoly je přiblížení způsobu, jakým byla data uložena v databázi a následně připravena k samotnému dolování z dat. Celý druhý oddíl je věnován přípravě dat pro samotné dolování z dat. Je zde popsán proces transformace atributů na proměnné kvantitativní. Poté následuje třetí podkapitola, která se podrobněji zabývá jednou z metod pro dolování dat a to *shlukováním*. Obsahem této části jsou již konkrétní algoritmy pro shlukování dat [29, 3] a také metoda *k-Means* využívaná v praktické části této práce. Kapitulu uzavírá stručný přehled vybraných programů pro data mining a podrobnější seznámení s nástrojem *RapidMiner*, který je v této práci využíván pro samotné dolování.

### Terminologie

Pojem data mining neboli česky dolování dat se začal ve vědeckých kruzích objevovat počátkem 90. let 20. století. První zmínka pochází z konferencí věnovaných umělé inteligenci (IJCAI'89<sup>1</sup> — mezinárodní konference konaná v Detroitu, AAAI'91<sup>2</sup> a AAAI'93 — americké konference v Kalifornii a Washingtonu, D. C) [2].

Tradiční metoda získání informací z dat je realizována jejich manuální analýzou a interpretací. V praxi ji například nalezneme v odvětví zdravotnictví, vědy, marketingu (efektivita reklamních kampaní, segmentace zákazníků) a dalších. Pro tyto a mnoho dalších disciplín je manuální zpracování příliš pomalé, drahé a vysoce subjektivní. Další důvod k přechodu na jiné metody je objemnost dat, která dramaticky vzrostla, a tudíž se manuální analýza stává zcela nepraktická. Databáze rychle rostou ve dvou následujících kategoriích:

1. počet záznamů neboli objektů v databázi
2. počet polí neboli atributů objektů v databázi

---

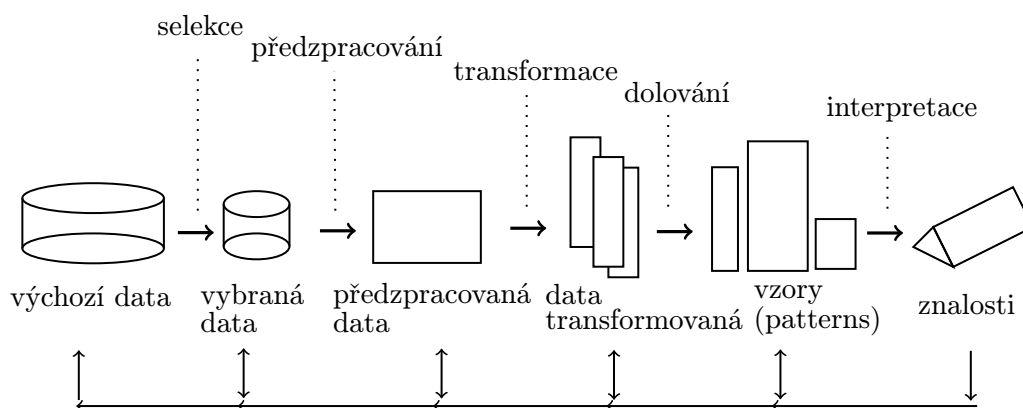
<sup>1</sup>International Joint Conference on Artificial Intelligence

<sup>2</sup>Association for the Advancement of Artificial Intelligence

Proces data mining je pouze jedna část z odvětví nazývané dobývání znalostí z databází neboli KDD<sup>3</sup> definované níže v definici 3.0.1. Vznik disciplíny KDD je důsledkem nepřehledného množství automaticky sbíraných dat, která je potřeba dále využívat. Podstatným znakem celého procesu je správnost reprezentace výsledků formou, která má k uživateli nejbližší. Jako příklad bude uvedena implikace ve tvaru rozhodovacích pravidel, asociační pravidla, rozhodovací stromy, shluky podobných dat a další. Základem KDD je praktická použitelnost metod. Očekává se zjištění nových skutečností namísto prezentování již známých informací.

**Definice 3.0.1** KDD je chápáno jako interaktivní a iterativní proces tvořený kroky selekce, předzpracováním, transformace, vlastního „dolování“ (data-mining viz 3.0.2) a interpretace [2].

Grafické znázornění definice 3.0.1 je popsáno schématem na obrázku 3.1, který prezentuje časový harmonogram v KDD. Schéma znázorňuje následnost jednotlivých procesů, které tvoří KDD. KDD je iterativní proces, z čehož vyplývá, že skutečnosti nalezené v předchozích částech zjednoduší a zpřesní vstupy pro následující fáze. Jakmile jsou znalosti získány, jsou prezentovány uživateli. Pro přesnost může být část procesu KDD ještě upravena. Tím budou získány „přesnější a vhodnější“ výsledky.



Obrázek 3.1: Proces dobývání znalostí z databází podle knihy autora Fayyad [4].

Vzhledem k obrázku 3.1, který prezentuje jednotlivé kroky získávání znalostí z databází, budou dále tyto procesy popsány. První část v KDD je tvořena výchozími daty, které slouží jako zdroj pro ostatní fáze. Samotný popis získávání těchto dat je popsán v předcházející kapitole. Procesu selekce dat je kladen za cíl vybrat co možná „nejúčinnější“ množinu dat a tím i zmenšit její celkový objem. V této části získávání znalostí se při vybírání dat bere ohled na to, jak se jednotlivá data vztahují ke konkrétnímu uživateli. Následující proces nazvaný transformace se zabývá převedením dat do vhodného formátu pro samotné dolování informací. Tato část je popsána v následující podkapitole, kde hlavní úlohu při transformaci je čas. Vybrané metody pro dolování dat jako je například shlukování a další, jsou taktéž v této práci popsány níže.

<sup>3</sup>Knowledge Discovery in Database

Získávání znalostí z databází je proces složený z několika kroků vedoucích od surových dat k formě nových poznatků. Iterativní proces je složený, jak je prezentováno v [5], z následujících kroků:

- **čištění dat** — fáze, ve které jsou nepodstatné údaje odstraněny z kolekce.
- **integrace dat** — kombinování heterogenních dat z několika zdrojů do společného jediného zdroje.
- **výběr dat** — rozhodování o relevantních datech.
- **transformace dat** — také známý jako konsolidace dat. Fáze, ve které jsou vybraná data transformována do formy vhodné pro dolování.
- **data mining** — zásadní krok, ve kterém jsou aplikovány vzory na data.
- **hodnocení modelů** — vzory dat zastupují získané znalosti.
- **prezentace znalostí** — konečná fáze, zjištěné poznatky jsou reprezentovány uživateli. Tento základní krok využívá vizualizační techniky, které pomáhají uživatelům porozumět a správně interpretovat získané výsledky.

Jak je uvedeno v [5], běžně jsou některé z těchto kroků kombinovány dohromady. Kroky čištění dat a integrace dat mohou být provedeny společně, tak jako to prezentuje schéma na obrázku 3.1.

V této podsekcí jsou ve stručnosti vysvětleny základní nejdůležitější pojmy dále v práci využívané.

Definice výrazu data mining se v odborné literatuře nachází několik. Zde uvedená je kombinací dvou „definic“ z [15].

**Definice 3.0.2** Data Mining je proces objevování znalostí, který používá různé analytické nástroje sloužící k odhalení dříve neznámých vztahů a informací z velmi rozsáhlých databází. Výsledkem je predikční model, který je podkladem pro rozhodování [15].

## 3.1 Transformace dat

Transformace dat tvoří třetí část z celkového procesu dobývání znalostí z databází. Než se data dostala do tohoto stavu, bylo na nich provedeno několik kroků, ve kterých byla upravována. V první fázi byla sbírána Jabber komunikace, která je popsána v první kapitole. Druhá fáze byla zaměřena na zúžení výsledné množiny, a proto byla vybrána jen relevantní data. I přes tyto kroky relační databáze obsahuje velké množství dat, která se nenachází ve stavu, aby mohla být použita jako zdroj pro data mining. Jak bude popsáno v následující podkapitole, většina metod pro dolování dat pracuje pouze s daty, která obsahují kvantitativní proměnné. Za tímto účelem je potřeba všechny atributy tabulky z databáze, které mají být nadále používány, převést na měřitelné hodnoty.

V této práci se bude pracovat s atributy nesoucí informace o jak aktuálních, tak minulých stavech uživatelů, kteří si přidali účet *jabInfo@jabim.com* do svého seznamu kontaktů. A tak byla jejich každá změna statusu uložena do databáze. Z důvodu nečíselné hodnoty stavů, jako je například *available*, *away* a další, je třeba provést jejich transformaci na kvantitativní hodnoty. Proces byl proveden pomocí bijektivního zobrazení. Kde zobrazení



je podle [8] funkce s definičním oborem  $S$  a oborem hodnot  $T$ , která je nazývána binární relace  $f \subseteq S \times T$ . V této relaci se nevyskytují dvě různé dvojice  $(s, t_1)$  a  $(s, t_2)$ , kde  $s \in S$  a  $t_1, t_2 \in T$ . Prvky  $t_1$  a  $t_2$  jsou různé. Z toho vyplývá, že každému prvku  $s$  z množiny  $S$  je přiřazen jednoznačně právě jeden prvek  $t \in T$ . Tuto definici je možné zapsat ve tvaru:

$$f : S \rightarrow T,$$

kde  $S$  a  $T$  jsou množiny  $(D_f, H_f)$ .

Konkrétní případ transformace z této práce tedy bude obsahovat množinu  $S$ , kde

$$S = \{Available, Chat, Away, DND, XA, Unavailable\}$$

a množina  $T$ , kde

$$T = \{120, 110, 90, 70, 50, 0\},$$

do které budou jednotlivé prvky z množiny  $S$  bijektivně zobrazeny. Kdy bijekce je zobrazení, které každému prvku z cílové množiny, konkrétně z množiny  $T$ , přiřazuje právě jeden prvek z množiny počáteční, tedy  $S$ .

Při výběru velikosti hodnoty, pro výslednou množinu  $T$ , bylo čerpáno z programu Gajim, který je multiplatformní klient s velkou podporou standardů a rozšiřujících protokolů. Hodnoty v množině  $T$  byly zvoleny tak, aby měly sestupné uspořádání, a aby bylo možné je dobře mezi sebou porovnávat. Jsou-li vybrány dvě hodnoty například available a unavailable, vzdálenost mezi nimi musí být větší než vzdálenost například u hodnot chat a away. Na druhou stranu prvky ze vstupní množiny  $S$  jsou striktně definovány podle Jabber standardu, který je popsán v RFC [23].

## Temporální data

Druhá podstatná transformace, pro kterou bylo čerpáno z [26], se tak jako první nezabývá transformováním dat textových na data, jejichž obsah by byl tvořen kvantitativními proměnnými. V této části jsou řídká temporální data transformována na hustá. Pro následné vyhodnocení a data mining je potřeba řádkům z tabulky presence přidat konečné časové razítko, které by vymezilo interval doby platnosti těchto dat.

Jak již bylo uvedeno, hlavním rozdílem mezi temporálními databázemi a ostatními je schopnost uchovávat časové údaje. Své uplatnění nachází v odvětvích, kde je potřeba zpracovávat stará a zároveň nová data, například v oblastech medicíny, finančnictví, monitorování a dalších. Jednotlivé záznamy, které jsou závislé na čase, jsou v databázích ukládány jako samotné body, tedy diskrétně. Přestože v reálném světě je většina těchto údajů z pohledu času spojitých.

K dalšímu popisu temporálních databází nyní budou charakterizovány tři důležité pojmy, pomocí nichž jsou databáze dále děleny. Prvním pojmem je *granualita*, která udává nejmenší časovou jednotku, kterou databáze rozlišují. Hodnoty, které může nabývat, jsou hodina, den, rok a další. Velikost granuality ovlivňuje velikost objemu dat, který je přímo úměrný s přesností záznamů. Dalším pojmem je *čas platnosti*, která reprezentuje období, kdy je daný fakt v modelovém světě pravdivý. Posledním termínem je *čas transakce*, která definuje přítomnost faktu v databázi a možnost jej získat. Čas transakce a čas platnosti jsou na sobě nezávislé a definují dvě rozdílné časové osy, kdy každá může disponovat s jinou granualitou. Souhrnný název pro výše uvedené tři pojmy, který se používá, je systém časových razítek. Při použití těchto systému lze následně tvořit dotazy zaměřené na různá časová



období jako je minulost, přítomnost a budoucnost. Tyto dotazy jsou velmi jednoduché a to díky rozšířenému jazyku TSQL, který vychází z klasické podoby dotazovacího jazyka SQL.

Temporální databáze jsou rozděleny podle systému časových razítek na tyto základní: *snímková*, *transakční*, *platného času*, *obojího času* (bitemporální). Entity z databáze, která je použita v této práci, je možné zařadit do kategorie *snímkových tabulek* (snapshot). K jejím hlavním rysům patří zaznamenávání stavu dat v jistém okamžiku. Čas transakce ani čas platnosti zde nejsou uplatněny.

Konkrétní příklad možné transformace je ukázán na části tabulky *presence* 3.1, která se ve stejném formátu nachází i v databázi, a modifikované tabulce *presence\_modify* 3.2. Tabulka 3.2 se od původní liší přidáním sloupce *dateEnd* (podbarven šedě), který s atri-

id	date	toj	presence
87365	2011-4-4 13:53:59	JabInfo@jabbim.cz	Available
87369	2011-4-4 15:43:07	JabInfo@jabbim.cz	Away
...	...	...	...

Tabulka 3.1: Ukázka tabulky *presence*.

butem *dateStart* vymezuje interval, kdy daná hodnota byla nebo je platná. Tato entita je pouze příkladem, jak by mohla daná transformace vypadat. V této práci se žádná nová tabulka nevytvářela a ani se nemodifikovala již vytvořená. Celá transformace je popsána v další části této podkapitoly.

id	dateStart	dateEnd	toj	presence
87365	2011-4-4 13:53:59	2011-4-4 15:43:07	JabInfo@jabbim.cz	Available
87369	2011-4-4 15:43:07	INF	JabInfo@jabbim.cz	Away
...	...	...	...	...

Tabulka 3.2: Ukázka modifikované tabulky *presence\_modify*.

Postup jednotlivých kroků při transformaci časových údajů z tabulky *presence* je zobrazen v následujícím výčtu. Tento zjednodušený popis algoritmu vyžaduje dva vstupní parametry. První je JabberID uživatele, jehož položky v databázi mají být transformovány. Druhá nutná položka je datum, která bude sloužit jako upřesňující vstupní interval.

Data z tabulky *presence* jsou transformována na vektor  $\vartheta$  o 288 dimenzích, kde z pohledu časového je jedna dimenze období vymezené 5 minutami. Den je rozdělen na úseky po 5 minutách ( $\delta = 300s$ ), kterých je 288. Tedy

$$\vartheta = (\vartheta_1, \vartheta_0, \dots, \vartheta_N),$$

kde  $N = 288$ .

1. Vstupním parametrem je datum, které vymezuje data pro transformaci. Toto datum je převedeno na dvě data (počátek a konec dne), která tvoří hraniční body v intervalu  $\iota$ .
2. Výběr dat z databáze, která splňují časové období definované intervalem  $\iota$  a uživatel  $ID$ . Uložení těchto dat do množiny  $\Gamma$ .
3. Pokud zadanému dotazu neodpovídá žádný řádek z tabulky, jsou data označena jako prázdná. Výstupní transformovaný vektor  $\vartheta$  je naplněn hodnotami reprezentujícími stav *Unavailable*. Algoritmus je **ukončen**.

4. Zjištění prvního statusu toho dne.
  - 4.1 Převod data o den dřívějšího na interval  $\iota_1$ , například  $\langle 2011-08-22\ 00:00:00, 2011-08-22\ 23:59:59 \rangle$ .
  - 4.2 Výběr dat vyhovujícím intervalu  $\iota_1$  a uživateli  $ID$  z databáze.
  - 4.3 Výběr posledního záznamu z množiny dat získaných z předešlého dotazu.
  - 4.4 Nalezení poslední presence a uložení její hodnoty do  $\lambda$ .
5. Výběr následujícího záznamu z množiny  $\Gamma$ .
6. Výpočet zda je časový interval mezi vybraným a následujícím záznamem větší jak  $\delta$ .
  - 6.a Časový interval je větší než interval  $\delta$ .
    - 6.1 Do výsledného vektoru  $\vartheta$  je ukládána hodnota presence daného záznamu.
    - 6.2 Opakuj předešlý bod **6.1** kolikrát je interval  $\delta$  menší než rozdíl mezi časy vybraného a následujícího záznamu.
  - 6.b Časový interval je menší než interval  $\delta$ .
    - 6.1 Jsou vybírány další záznamy z množiny  $\Gamma$ , dokud rozdíl mezi časy v sekundách není větší než interval  $\delta$ .
    - 6.2 Jednotlivé presence záznamů jsou ukládány do pomocného pole, transformované do kvantitativních hodnot, jak je uvedeno v úvodu této podkapitoly.
    - 6.3 Každý prvek pole je vynásoben počtem sekund, zastoupených v daném intervalu.
    - 6.4 Výběr největšího prvku z pomocného pole a uložení jej do výsledného vektoru  $\vartheta$ .
7. Celý proces je opakován od bodu **5**, dokud množina  $\Gamma$  není prázdná.

## 3.2 Shlukování

Základ metod dolování dat je založen na statistice, posledních poznatcích z umělé inteligence či strojového učení. Hlavní cíl těchto netriviálních metod je společný — snaha zjištěné výsledky prezentovat srozumitelnou formou. Pro většinu používaných metod je společná vlastnost předpoklad, že objekty popsané pomocí podobných charakteristik patří do stejné skupiny (učení na základě podobnosti similarity-based learning). Objekty obsahující atributy lze převést na body v  $n$ -rozměrném prostoru, kde  $n$  reprezentuje počet atributů. Vychází se z představy podobnosti bodů tvořící určité shluky v prostoru.

Další rozdíly mezi metodami, které byly prezentovány v [2], spočívají v:

- schopnosti reprezentace shluků (např. otázka lineární separability)
- srozumitelnosti nalezených znalostí pro uživatele (symbolické vs. subsymbolické metody)
- efektivnosti znovupoužití nalezených znalostí
- vhodnosti typů dat

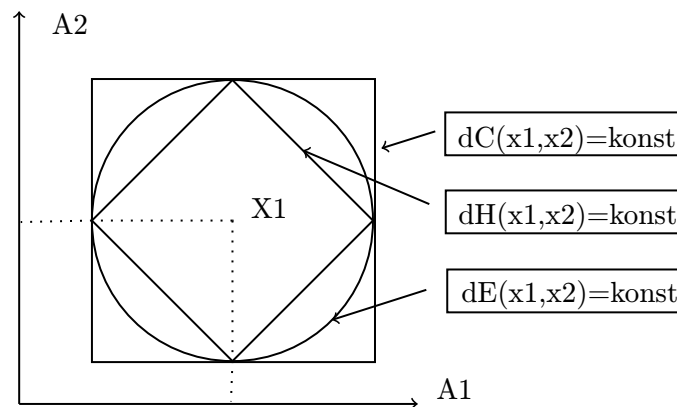
- a další ...

Problémy, které data mining řeší, se rozdělují do několika skupin. Do výčtu vybraných z nich, které budou následně rozebrány, patří *asociační pravidla*, *klasifikace*, *modely*, *predikce* a *shlukování*.

Následující část se bude zabývat shlukováním, které je rozděleno na několik metod shlukové analýzy podle [5]. U každé z nich jsou popsány její základní vlastnosti a uvedeny algoritmy relevantní k práci. Poslední metoda *metoda rozkladu* je rozebrána podrobněji z důvodu jejího praktického využití v této práci.

Shlukování je zaměřeno na dělení objektů do předem neznámých skupin. Proces dělení probíhá na základě specifikace objektů a jejich odlišnosti od ostatních shluků.

Většina níže popsaných metod a algoritmů je založena na výpočtu vzdáleností mezi objekty. Tato vzdálenost lze vyjádřit různými mírami, podle knihy [2] například pomocí *Hammingovy vzdálenosti* (dH), *Euklidovské vzdálenosti* (dE) a *Čebyševovy vzdálenosti* (dC). Rozdíl mezi těmito typy určující vzdálenosti graficky vyjadřuje obrázek 3.2. Kde  $X_1$  je střed, od něhož jsou jednotlivými obrazy znázorněny dané vzdálenosti. Konkrétně pomyslné body umístěné po obvodu kruhu jsou všechny stejně vzdáleny od středu  $X_1$ . Tato vzdálenost je označena jako Euklidovská. Další 2D těleso čtverec, který je vodorovný s osami  $A_1$  a  $A_2$  prezentuje Čebyševovu vzdálenost. Po obvodu posledního obrazce, čtverce otočeného o  $45^\circ$  podle osy  $A_1$ , jsou všechny pomyslné body stejně vzdáleny od bodu  $X_1$  Hammingovou vzdáleností.



Obrázek 3.2: Srovnání výpočtu vzdáleností od bodu  $x_1$  [2].

Jako první jsou zde rozebrány *metody založené na modelu*, které se pokouší přiřadit data k určitému matematickému modelu na základě společných optimalizovaných vlastností. Většina procesů je založena na předpokladu, že jsou data generována pomocí standardních statistik. Mezi zástupné metody této shlukovací analýzy se řadí Expectation-Maximization (EM) a Self Organizing Oscillator Network, dále jen SOON. Algoritmus SOON je založen na neuronové síti. Je to metoda vycházející z algoritmu SOM<sup>4</sup> [29]. Metoda EM je rozšířením algoritmu *k-means*, který bude podrobně rozebrán v následující části.

Hlavní princip *metody hierarchického shlukování* je založen na tvorbě stromové hierarchie shluků, která je známá pod názvem *dendrogram*. Hierarchické metody, podle [5], mohou být rozděleny do dvou skupin a to na základě principu, kterým jsou dendrogramy vytvářeny.

<sup>4</sup>Self-Organizing Map

První možnost je *aglomerativní přístup*, který shlukuje menší shluky, kdy výsledkem je jen jeden. Druhý přístup, *divizní*, je založen na opačném předpokladu. Na počátku je tedy jeden velký shluk, který je postupně rozdělován, dokud není počet shluků roven počtu objektů [29]. Mezi zástupce této metody například patří algoritmus AGNES<sup>5</sup>.

K dalším metodám patří *metody založené na mřížce*, které kvantují datový prostor do konečného počtu pravoúhlých buněk. Tyto buňky jsou uspořádány do víceúrovňové mřížkové struktury. Zmíněná struktura tvoří základ pro shlukové operace. Hlavní výhodou tohoto přístupu je rychlost zpracování, které většinou nebere ohled na počet datových objektů. Čas zpracování závisí pouze na počtu buněk v každé dimenzi kvantovaného prostoru. Mezi zástupce metod založených na mřížce patří metoda STING — STatistical INformation Grid, který pracuje se statickými informacemi uloženými v buňkách mřížky. Algoritmus je rozdělen do dvou částí. První si klade za cíl rekurzivně rozdělit datový prostor na pravoúhlé buňky. Druhá fáze testuje spojitost mezi sousedy relevantních buněk [29]. Mezi další metody založené na mřížce patří WaveCluster<sup>6</sup> využívající vlnkové transformace k rozdělení prostoru dat. Tato transformace zdůrazňuje shluky v prostoru a objekty jim vzdálené potlačuje [5].

*Metody založené na hustotě* vychází z  $m$ -rozměrného prostoru, ve kterém jsou zobrazeny objekty ve formě bodů. Místa v prostoru s větší koncentrací objektů ve srovnání s ostatními oblastmi jsou nazývány shluky. Výchozí předpoklad je existence okolí jednotlivých bodů (sousedství). Jedna z charakteristik metod založených na hustotě je schopnost vyřadit se s vzdálenými hodnotami, označovanými jako šum [29]. Jako příklad je uvedena metoda DBSCAN<sup>7</sup>, která je založena na hustotě objektů v prostoru. U jednotlivých objektů je zkoumáno jejich okolí. Algoritmus je ovlivňován dvěma parametry, velikostí shluku  $\varepsilon$  a minimálním počtem objektů v daném shluku  $MinPts$ , které spolu úzce souvisí (viz [5]). Bod splňující obě podmínky je označen za jádro. Za pomoci jader je rozšiřována množina objektů spojených na základě hustoty. Obsahuje-li jádro  $x_1$  ve svém okolí další centrum  $x_2$ , znamená to, že  $x_1$  je přímo dosažitelné z  $x_2$ . Tímto způsobem jsou vytvářeny výsledné shluky. V opačném případě body, které nesplňují dvě zmíněné podmínky, jsou označeny jako šum.

Všechny zde doposud zmiňované metody poskytují dobré výsledky pouze s malým počtem dimenzí tak, jak je to popsáno v [8]. S narůstajícím počtem atributů roste počet nerelevantních dimenzí určených pro shlukování. S tímto také přibývá zvětšená produkce zašumění a znesnadnění nalezení relevantních shluků. Data jsou roztroušena do mnoha dimenzí a tím odpadá možnost použití vzdálenostních funkcí. Zmíněné problémy shlukování velkých dat řeší dvě techniky, které se nazývají: *metoda transformace rysů* a *metoda výběru atributů*. Pro efektivní shlukování je možné použít například algoritmus CLIQUE<sup>8</sup>.

## Metody rozkladu

Metody rozkladu rozdělují datové prvky do několika podmnožin nazývané shluky. Počet shluků musí být znám před zahájením samotného procesu. Přiřazení do konkrétních tříd je podle [29] jednoznačné nebo probíhá na základě míry příslušnosti objektů do shluků. Pro velký počet objektů, se kterými se pracuje, jsou využívány různé iterační optimalizace.

Hlavním zástupcem u uvedených metod je algoritmus  $k$ -means, který je popsán níže.

---

<sup>5</sup>AGglomerative Nesting

<sup>6</sup>Clustering Using Wavelet Transformation

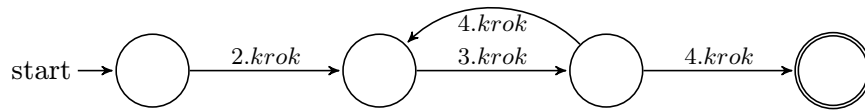
<sup>7</sup>Density-Based Spatial Clustering of Applications with Noise

<sup>8</sup>CLustering In QUEst

Tvoří základ pro většinu metod shlukování nejen pro metody rozkladu. K dalším metodám se řadí  $k$ -medoidů,  $k$ -modů,  $k$ -histogramů, fuzzy shluková analýza a další.

### $k$ -means

Shlukování pomocí algoritmu  $k$ -means je používáno pro data obsahující kvantitativní proměnné a pro data, která nejsou příliš zašumělá. Základní proces je tvořen iterativním rozdělováním objektů do tříd na základě vzdáleností od jejich středů. Střed, neboli centroid shluku, je vektor, jehož vzdálenost od součtu vzdáleností objektů v této třídě je minimální. Celý tento proces je prezentován na obrázku 3.3 pomocí jednoduchého schématu konečného automatu. Jednotlivé kroky konečného automatu odpovídají krokům  $k$ -means zobrazeného pomocí algoritmu 3.1. Pro výpočet vzdáleností mezi objekty samotnými nebo mezi objekty



Obrázek 3.3: Algoritmus  $k$ -means zobrazený pomocí konečného automatu.

a středem je použita euklidovská vzdálenost<sup>9</sup>, která je vyobrazena na obrázku 3.2.

K hlavním výhodám algoritmu  $k$ -means patří jeho relativní efektivnost. Složitost algoritmu je  $O(TKN)$ , kde  $N$  je počet objektů,  $K$  je počet shluků a  $T$  je počet iterací. Obvykle platí, že počet objektů je mnohem větší než počet iterací i shluků. Na druhou stranu má i řadu nevýhod, kvůli kterým je často různými způsoby modifikován ( $k$ -medoids,  $k$ -medians). K hlavním „nedostatkům“ patří předem nutná znalost počtu shluků (tříd)  $K$ , do kterých budou objekty zařazeny. Druhý často se vyskytující problém je samotné ukončení algoritmu, které nastane u nalezení lokálního optima namísto optima globálního. Tato nepřesnost vzniká nevhodně zvoleným rozmístěním počátečních středů. Původní nemodifikovaná verze algoritmu nedefinuje, jak se má postupovat, jsou-li nalezeny prázdné shluky.

$K$ -menas je algoritmus, kterým jsou přiřazovány objekty (vektory)  $x_n$ , kde  $n = 1, \dots, N$ , do  $S_k$ , kde  $k = 1, \dots, K$ , shluků. V prvním kroku jsou určeny počáteční středy tříd, do kterých se budou objekty shlukovat. Určení počátečních centroidů  $c_k$  probíhá například náhodným výběrem  $K$  objektů nebo  $K$  prvních objektů souboru. Druhým krokem jsou zkoumány jednotlivé vzdálenosti objektů  $x_n$  od počátečních středů  $c_j$  pomocí euklidovské vzdálenosti. Na základě nejmenší zjištěné vzdálenosti mezi objektem a centroidem je objekt zařazen do shluku, kterému náleží právě tento střed. Ve třetím kroku, tak jako u kroku prvního, jsou hledány nové středy shluků. Nyní však již nejsou zvoleny náhodně, ale spočítány. Jsou vypočítány na základě průměrných jednotlivých hodnot objektů a uložen jako  $m$ -rozměrný vektor. Čtvrtým krokem se algoritmus dostává do konečné fáze, kdy mohou nastat dva možné případy. Nově nalezené středy nejsou příliš vzdáleny od předchozích centroidů a proto je algoritmus ukončen. Druhá častěji se vyskytující možnost iterativně provádí algoritmus od druhého kroku, dokud neplatí první možnost nebo dokud se objekty nepřestanou přemisťovat úplně. Při popisu tohoto algoritmu bylo čerpáno z [29, 2]. Níže zobrazený algoritmus 3.1 prezentuje krok po kroku metodu  $k$ -means.

<sup>9</sup>mean = střed, centroid je vektor průměrů

- 
1. náhodně zvol rozklad do  $K$  shluků
  2. urči centroidy pro všechny shluky v aktuálním rozkladu
  3. pro každý příklad  $x$ 
    - 3.1 urči vzdálenosti  $d(x, c_k)$ ,  $k = 1, \dots, K$ , kde  $c_k$  je centorid  $k$ -tého shluku
    - 3.2 nechť  $d(x, c_l) = \min_k d(x, c_k)$
    - 3.3 není-li  $x$  součástí shluku  $l$  ( $k$ -jehož centoridu  $c_l$  má nejblíže), přesuň  $x$  do shluku  $l$
  4. došlo-li k nějakému přesunu, potom jdi na 2, jinak konec
- 

Algoritmus 3.1: Metoda  $k$ -means byla převzata z [2].

Díky jednoduchosti a relativní rychlosti je metoda  $k$ -means stále výrazně využívána. Uplatnění nachází v široké škále oblastí jako je například biologie nebo počítačová grafika. Vzhledem k enormnímu počtu možného uspořádání nejsou výsledky vždy přesné, ale často pouze přibližné.

### 3.3 Programy

V současné době na programovém trhu existuje mnoho systému, které jsou zaměřeny na data mining. Mezi nejrozšířenější a nejdostupnější nástroje patří Weka a RapidMiner. K těmto nástrojům je také možné zařadit program FIT-miner vyvíjený na fakultě informačních technologií v Brně. V této práci bylo pro samotné dolování z dat využito programu RapidMiner, který dostal přednost před ostatními. Z pohledu nástroje FIT-miner, který ve své základní části podporuje z databází pouze Oracle, se RapidMiner jevil jako vhodnější. Kompatibilitu pro databáze typu PostgreSQL již měl zabudovanou a tak nebylo potřeba vyvíjet žádné doplňující moduly, jak by to bylo u FIT-mineru. V případě nástroje Weka, RapidMiner působil propracovanějším dojmem a také nabízí lepší grafické zobrazení vyhodnocených výsledků.

Dalším velmi rozšířeným a často používaným nástrojem je jazyk  $R$ .  $R$  vychází z jazyka  $S$ , který ale není jako jazyk  $R$  volně šiřitelný. Statistický a grafický nástroj  $R$  je tedy volně dostupným jazykem a prostředím, které je ovládáno pouze z příkazové řádky. Pro jednodušší práci jej lze rozšířit o grafické rozhraní jako je RkWard nebo R Commander. Samotnou aplikaci lze rozšířit o mnoho statistických doplňků, které jsou taktéž zdarma.

Mnoho programů pro dolování dat je založena na přístupu vizuálního programování. Jedná se o proces, při kterém je uživatelem za pomoci grafických prostředků navržen algoritmus a další postup práce. Jako příklad lze uvést poloprofesionální aplikaci *Orange*, u které jsou nejdůležitější části psány pomocí C++ a rozšíření lze implementovat v jazyce Python.

Na vybraných technicky zaměřených vysokých školách existují skupiny, které se zabývají výzkumem a vývojem nástrojů pro data mining. Jako příklad lze uvést již dříve zmiňovaný FIT-miner z fakulty informačních technologií v Brně nebo také projekt LISp-Miner z Vysoké školy ekonomické v Praze. LISp-Miner je otevřený akademický systém určený pro výuku a výzkum metod pro dobývání znalostí z databází.

#### RapidMiner

RapidMiner je, tak jak je popsán na oficiálních stránkách produktu [28], celosvětově nejpožívanější open-source systém pro dolování dat. Je možné jej používat jako samotnou aplikaci nebo jej začlenit jako komponentu do vlastních výrobků v podobě knihovny pro jazyk Java. Pro zájemce je nabízen také ve verzích pro firmy, které jsou rozdílné v poplatcích,

podpoře pro zákazníka, záruce a dalších balíčků služeb zajišťujících celkovou komplexnost a spolehlivost produktu.

Jako již většina podobných aplikací je RapidMiner v současné době implementován v jazyce Java, díky které nabízí flexibilní nejen grafické prostředí. K vybraným základním rysům tohoto nástroje, tak jak jsou prezentovány firmou *Rapid-i*, patří: výkonné, přesto intuitivní grafické uživatelské rozhraní pro návrh procesů, jednoduché řešení pro transformaci dat, kontrola výsledků již při samotném návrhu a další. Nástroj RapidMiner podporuje širokou škálu metod a algoritmů pro data mining. Mnoho algoritmů je implementováno přímo v aplikaci, ale také je použito metod z konkurenčního softwaru Weka. V základní verzi určené pro veřejnost je k nalezení přes 100 procesů k modelování. Jsou zde zastoupeny jak metody klasifikační a asociační, tak i metody shlukovací, z nichž lze jmenovat například DBSCAN,  $k$ -medoids a hlavně  $k$ -means.

K dalším schopnostem RapidMineru je možnost spuštění jeho samotného pomocí grafického rozhraní nebo z příkazové řádky. Jak již bylo uvedeno dříve, je také možné jej použít jako knihovnu v jazyce Java. V této práci jsou použity první dvě možnosti. Pomocí grafického prostředí byl vytvořen experiment, otestována jeho funkčnost a následně pro jednotlivá shlukování použita šablona procesu, která byla volána z příkazové řádky. Tato možnost je k dispozici díky tomu, že jsou projekty v programu RapidMiner ukládány do čitelné a strukturované formy za pomoci značkovacího jazyka XML.



## Kapitola 4

# Implementace

Obsahem čtvrté kapitoly je popis praktické části této práce. Jsou zde charakterizovány jednotlivé prvky, které byly použity jak pro získání dat, tak pro jejich následné uložení. V první části je prezentována struktura architektury této práce. Její grafické znázornění je ukázáno na obrázku 4.1.

Cílem této práce je dolování dat z Jabberu. Jak již bylo dříve napsáno, Jabber je síťová služba pracující v reálném čase díky níž mohou její uživatelé komunikovat, informovat nebo sdílet svůj status s jinými uživateli. Celá tato vzájemná komunikace skrývá rozsáhlé množství informací o klientech dané sítě. Všechna tato navzájem vyměněná nebo poskytnutá data následně poslouží jako zdroj samotnému dolování. Pro jejich uskladnění je využita databáze, jejichž strukturální návrh prezentuje obrázek 4.2.

Třetí část této kapitoly je zaměřena na Jabber klienta neboli robota, který je implementován v jazyce C++ pouze s konzolovým rozhraním. Robot v této práci hraje roli pasivního uživatele, který informace pouze přijímá. Ve vybraných případech dokáže uživatele ze svého seznamu kontaktů vyzvat k zaslání odpovědi s informacemi na odpověď, o kterou žádal. Struktura samotného robota je popsána níže a reprezentována obrázkem 4.3. V části o Jabber robotovi jsou také uvedeny informace o knihovně *gloox*, kterou jsou zprostředkovány všechny náležitosti Jabber komunikace.

Poslední podkapitola této části je věnována případnému pokračování této práce.

### 4.1 Architektura

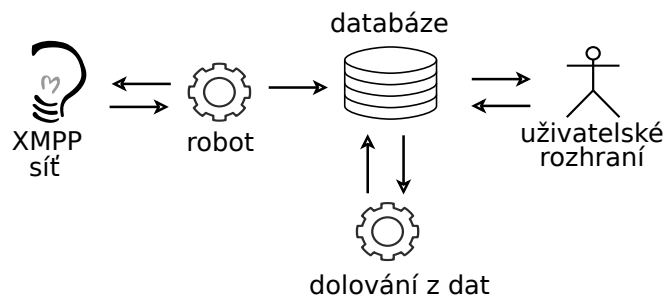
Již ze samotného zadání a názvu této technické zprávy je zřejmé, že je třeba celý proces rozčlenit na menší elementy. Vhodnou dekompozicí vzniklo pět jednotlivých ucelených prvků, které jsou prezentovány schématem na obrázku 4.1. Konkrétně to jsou, zleva: *XMPP server*, *robot*, *databáze*, *data mining* a *uživatelské rozhraní*. Vzájemná výměna informací mezi jednotlivými částmi je zobrazena pomocí šipek, které určují směr komunikace.

V levé části obrázku 4.1 je blok XMPP síť, který prezentuje zdroj dat. Fungování XMPP sítě je detailně popsáno ve druhé kapitole. Vztah a vzájemná komunikace s robotem probíhá pomocí internetové sítě, kdy XMPP síť souhrnně reprezentuje jednotlivé prvky, jako jsou Jabber servery, uživatelé a jejich klienty. Jak je patrné, výměna informací mezi těmito částmi probíhá obousměrně. Druhý element schématu *robot*, je charakterizován v podkapitole níže, kde je podrobně popsán návrh jeho struktury. Jsou zde zdůrazněny jeho základní vlastnosti a schopnosti. Dalším blokem je *databáze*, která reprezentuje datové úložiště. Do něhož jsou za pomoci jednosměrného kanálu ukládána data, která jsou získávána z toku



informací proudícími mezi robotem a XMPP sítí. Přehled vybraných jednotlivých tabulek a jejich atributů je možné nalézt v následující kapitole, která se zabývá návrhem databáze.

Dalším oddílem, prezentovaným na obrázku 4.1 pod názvem *data mining skript*, je skript, který provádí samotné dolování z dat. Jako první jsou vybraná data z databáze transformována do vhodného formátu pro dolování z dat. Tato část je podrobně popsána ve třetí kapitole, v oddílu zabývajícím se transformací dat. Přeformátovaná data jsou předána programu RapidMineru, kterým za pomoci algoritmu *k*-means je prováděn data mining.



Obrázek 4.1: Struktura architektury bakalářské práce.

Poslední částí je *úživatel'ské rozhraní*, které je implementováno pomocí webových služeb. Jako příklad lze uvést službu, která zobrazuje status uživatele na webové stránce. Uživateli pouze stačí vlastnit Jabber účet a mít přidáného tohoto robota do seznamu kontaktů. Od každého uživatele, jež má tohoto robota přidáného do svého seznamů kontaktů, je sbírána veškerá síťová aktivita. Možnosti, které mu nabízí uživatelské rozhraní v podobě webové prezentace, tudíž záleží pouze na datech, které sám do sítě rozesílá. Jsou-li podporována všechna zde implementovaná rozšíření, která jsou uvedena v kapitole druhé, je možné využít jejich služeb prostřednictvím internetových stránek. Jako další příklad lze uvést zobrazení aktuálního umístění uživatele na mapách od firmy google, dle informací poskytnutých prostřednictvím dokumentu User Geoloc [6].

## 4.2 Návrh databáze

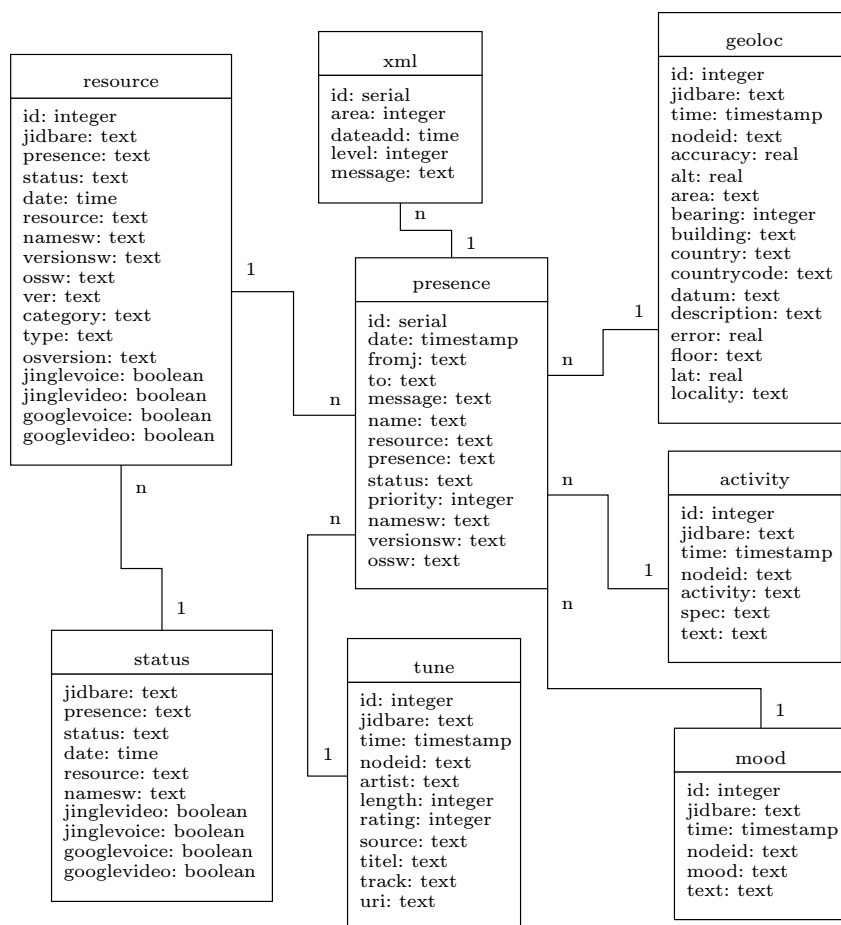
Data, která jsou sbírána robotem, jsou ukládána do objektově-relační databáze PostgreSQL. PostgreSQL neboli také Postgres byl použit ve verzi 8.4.7 a je provozován na operačním systému Ubuntu.

Struktura databáze, do které je ukládána veškerá komunikace Jabber robota, využívá relační model. Obrázkem 4.2 jsou prezentovány nejdůležitější části databáze. Celou strukturu návrhu je možné nalézt v příloze D. V jednotlivých částech návrhu databáze je počítáno s druhotným využitím obsahu, které bude popsáno v dalších oddílech této práce.

V době návrhu databáze nebylo zcela zřejmé, která data budou následně analyzována. Z tohoto důvodu se struktura databáze snaží zachytit všechna „důležitá“ data. Za tímto účelem je v návrhu databáze obsažena tabulka *xml*, která je nositelem obsahu jak všech přijatých, tak i odeslaných zpráv. Tabulky *debug*, *level* a *logarea*, které v zúženém návrhu databáze, uvedeném na obrázku 4.2, nejsou zobrazeny, jsou určeny pouze jako doplňkové informace k typu zprávy. Tabulku *xml* je možné nahradit jednotlivými dalšími tabulkami, které jsou zaměřeny na konkrétní data. Příkladem entity, která již obsahuje konkrétní data bez XML prvků, je tabulka *message*, která je taktéž vyobrazena v příloze C. Jejím obsahem

jsou zprávy vzniklé při Jabber komunikaci mezi uživatelem a robotem, jehož schopnosti budou popsány níže.

Tabulka *presence* z pohledu XMPP standardu prezentuje jeden ze tří základních částí stanzy, element *presence*. Základním cílem této tabulky je shromažďování uživatelských statusů. Jak již bylo zmíněno ve třetí kapitole, v části která se zabývá transformací, atribut *presence* obsahuje hodnoty typu text. Příklad všech možných hodnot, kterých tento atribut může nabývat je prezentován v příloze C v části zabývající se elementem *presence*. K dalším atributům této tabulky patří *message*, který je reprezentován textovým řetězcem a rozšiřuje informace o stavu uživatele. V této části je často obsažen text, který je do statusu přidán automaticky IM klientem. Transformovaný obsah této entity tvoří důležitou část při získávání znalostí a to v podobě dat, ze kterých budou „dolovány“ informace.



Obrázek 4.2: Vybraná část struktury databáze.

Většina rozšíření standardu XMPP, která jsou popsána ve druhé kapitole, jsou prezentována pomocí pěti tabulek. Konkrétně to jsou tabulky *geoloc*, *tune*, *mood*, *activity* a *vcard*. Pro obsah a názvy atributů všech pěti tabulek s rozšířeními se staly vzory dokumenty jednotlivých XEP, které je definují. Tabulky kromě těchto atributů obsahují také položku *jidbare*, která, jak již název naznačuje, obsahuje pouze čisté JabberID bez resources. Tato skutečnost vyplývá již ze samotného návrhu XEP protokolů. V tabulce *vcard* jsou také položky, jejichž obsah je tvořen fotografiemi. Obrázky jsou zde uloženy ve stavu, tak jak

jsou přenášeny po síti, tedy pomocí datového formátu Base64. Base64 je algoritmus, který převádí binární data na řetězec, který je tvořen pouze znaky ASCII tabulky.

S posledním rozšířením, nazvaném *Software Version*, je spjata tabulka resource, jejíž několik sloupců tvoří informace o IM klientovi a operačním systému, na kterém běží. Tyto základní údaje jsou definovány v protokolu [17], kde je také možné čerpat bližší informace. Pokročilejší atributy, jako je *type* a *osversion* nacházejí svůj podklad v XEP dokumentu [7], taktéž jako atribut *ver*. V neposlední řadě se zde nachází zmínka o podpoře audia a videa, ať už v podobě *jingle* nebo *google*. Přesto tyto výše uvedené údaje nejsou hlavním důvodem existence této tabulky. Její primární cíl je uchovávat informace o připojených uživateliích a všech jejich resources.

Entita resources tvoří základ pro tabulku *status*. Její obsah je automaticky generován z obsahu tabulky popsané výše. Počet řádků odpovídá počtu uživatelů, kteří jsou v seznamu kontaktů tohoto robota. Primárním cílem je informování o aktuálním stavu uživatele. Je tedy poskytován stav, ve kterém se uživatel nachází s přihlédnutím na prioritu a resources. Řádek entity obsahuje status a dostupnost uživatele a jeho resources s největší přihlášenou prioritou.

### 4.3 Návrh robota

V této části bude popsána aplikace, která zajišťuje real-time komunikaci s ostatními entitami Jabber sítě. Konzolový robot, který je vyvíjen pro operační systém Linux, také vhodným způsobem pracuje s databází, kam jsou ukládána jednotlivá získaná data. Jabber jádro síťové komunikace [23] je implementováno pomocí volně dostupné knihovny *gloox*. V porovnání s jinými knihovnami disponuje lepší podporou a dokumentací. Gloox plně implementuje standart XMPP Core [23] a z větší části i standard XMPP IM [24]. Dodatečně je podporováno kolem třiceti XEP standardů mezi něž například patří vcard-temp [16].

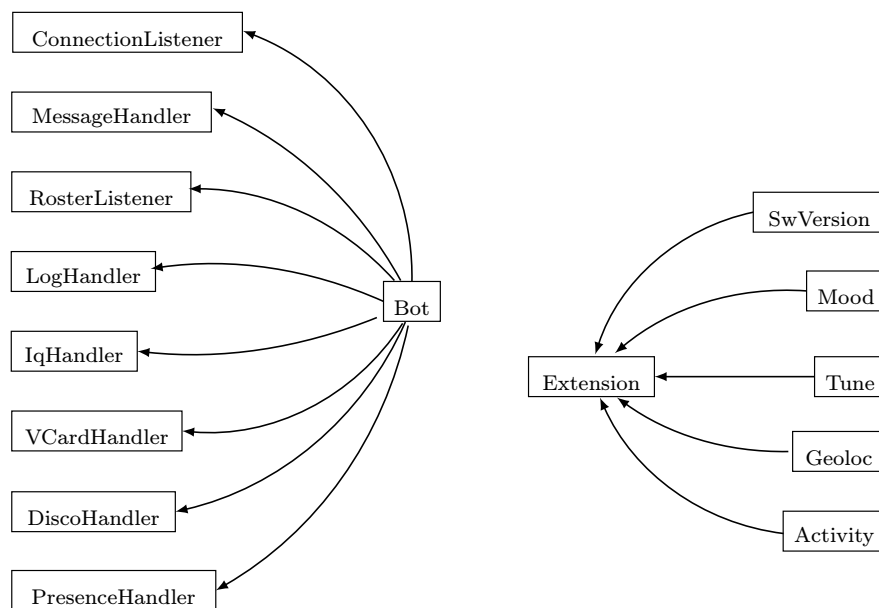
Při implementaci robota byly využity základní prvky objektově-orientovaného programování. Každá třída zde zastupuje jednoznačně oddělitelné elementy robota. Ať už se jedná o blok komunikující s databází, různá rozšíření nebo o jádro robota. Jako příklad lze uvést třídy, které reprezentují jednotlivá pro tuto práci vhodná rozšíření, jejichž základ je v podobě XEP dokumentů.

Struktura návrhu tříd robota, která je v základní části prezentována obrázkem 4.3, je velmi podobná struktuře návrhu databáze, obrázek 4.2. Mnoho atributů z robota a z databáze jsou ve vztahu jedna ku jedné. Samotná komunikace a řízení mezi těmito entitami je zprostředkováno pomocí knihovny *libpq*. V souboru *connect.h* jsou, v podobě předdefinovaných konstant, uvedeny všechny údaje nutné k navázání spojení ať už s databází nebo s komunikačním Jabber serverem.

Komunikace s databází se nachází v samostatné třídě, která zajišťuje provázanost mezi ní a daty. Pomocí funkcí jsou data, získaná z Jabber komunikace, přenášena a ukládána do jednotlivých tabulek databáze. Při zahájení činnosti robota je zajištěno navázání připojení k databázi a provedení nutných inicializačních kroků. Jako je tvorba nových tabulek, při prvním zapnutí aplikace, nebo kontrola existence těchto entit.

Pro rozšíření podle dokumentů XEP popsaných v kapitole druhé, v části zabývající se implementovanými rozšířeními, je základ nadtřída *Extension*. Tato třída obsahuje základní funkce a parametry pro všechna rozšíření. Třídy, které z *Extension* dědí a jsou prezentovány na návrhu robota uvedeném na obrázku 4.3, jsou následující: *Geoloc*, *Tune*, *Mood*, *Activity* a *SwVersion*. Poslední doposud nezmiňované rozšíření v podobě *vcard*, je využito z knihovny *gloox*, která jej implementuje.

Základní část, kterou je možné považovat za jádro aplikace, je třída *Bot*. Je implementována pomocí tzv. „handler“ tříd z knihovny gloox, kterými je zajišťován přístup k jednotlivým zprávám. Neboli jejich prostřednictvím je získán přístup k samotným elementům stanzy, jako je message, iq, a presence. Konkrétní seznam tříd, ze kterých je děděno, zobrazuje návrh robota na obrázku 4.3. Jako příklad lze uvést třídu, z prostoru jmen gloox, *VcardHandler*, díky níž je možné získat a zpracovávat vcard uživatelů, kteří jsou v seznamu kontaktů robota. Samotné přidávání uživatelů do seznamu kontaktů je prováděno automaticky. Pouze na straně klienta je nutné požádat robota o autorizaci, která však proběhne automaticky. Jednotlivé kontakty nejsou žádným způsobem tříděny do skupin, tudíž existuje pouze jedna.



Obrázek 4.3: Vybraná část struktury robota.

K dalším schopnostem robota, tedy kromě sbírání dat a jejich následného ukládání do databáze, patří reagovat na vybrané zprávy. Tyto zprávy jsou přijímány pouze od uživatelů, kteří byli začleněni do seznamu kontaktů tohoto robota. Základní příkaz, který nesmí chybět u žádné aplikace, je získání nápovědy. V první řadě informuje o verzi robota a jeho základních vlastnostech. Dále je odeslán stručný, ale výstižný seznam příkazů i s jejich vysvětlením. Při zaslání zprávy s neznámým příkazem robotovi, je uživateli vrácen stručný seznam kompatibilních příkazů. Přehled těchto příkazů se nachází v příloze E.

## 4.4 Pokračování práce

Obsah této podkapitoly je zaměřen na možnosti rozšíření robota sbírajícího data. Také jsou zde popsány možné změny, které se týkají uživatelského rozhraní. V závěru jsou shrnuta další rozšíření.

Většina služeb, které jsou dostupné na internetu a jsou podporované různými servery, poskytuje pouze informace o aktuálním stavu klienta. Pomocí výše uvedených standardů by bylo možné připojit k základní webové prezentaci statusu i informace rozšířené. Také zapojení standardu *Entity Capabilities* [7], který je popsán na konci druhé kapitoly, by mohlo

přinést nové pokročilejší, užitečné informace. Konkrétně je to samotná podpora IM aplikací v oblasti dalších XEP dokumentů. Nezasvěcený uživatel by tímto velice jednoduchým rozšířením dokázal zjistit seznam služeb podporovaných jeho IM programem. Ať už by se jednalo o seznam funkcí, které jsou pouze přijímány, nebo výčet rozšíření, k jejichž odběru je zaregistrován. Jako příklad lze uvést zjištění schopnosti, ať už přijímání nebo odesílání informací o přehrávané hudbě popsané dokumentem *User Tune* [18].

Další oblast, ve které lze nalézt možnost rozšíření, se také zabývá webovou prezentací. Nyní však jde již konkrétně o robota vytvořeného k účelu této práce. Informace potřebné k zobrazení aktuálního stavu klienta na internetové stránce, jsou uloženy v databázi v tabulce status. Tato tabulka je automaticky generována z obsahu jiné entity, která obsahuje údaje o všech připojeních daného uživatele. Vše pracuje správně až do té chvíle, kdy uživatel svou IM aplikaci neukončí korektně. Nastává situace, kdy na server není zaslána zpráva, která informuje o změně statusu. Díky tomuto neočekávanému ukončení klienta je uživatel prezentován, jako by byl připojen i když je tomu naopak. Zdrojová tabulka se v této chvíli stává neaktuální. Navrhované řešení této situace využívá rozšíření *XMPP ping* [20], díky němuž by se mělo dát zjistit, zda je uživatel stále dostupný nebo již nikoliv.

V oddílu robota by bylo možné pokračovat v části, která se zabývá příkazy zaslané ze strany uživatele. S tímto rozšířením úzce souvisí rozčlenění uživatelů do jednotlivých skupin. Každá takováto část seznamu kontaktů by vlastnila různá práva, která by se stala základem pro poskytování konkrétních služeb. Zařazení do těchto skupin by mohlo probíhat při žádosti o počáteční autorizaci, která by obsahovala zprávu s přesně definovaným řetězcem. Funkce, která je již uvedena výše a to v podobě rozšíření webové prezentace, by také mohla být poskytována prostřednictvím robota. Konkrétně se jedná o zjištění podporovaných ať již přijímaných nebo odesílaných rozšíření. Uživatel by pouze poslal robotovi striktně předdefinovanou zprávu a jako odpověď by obdržel seznam podporovaných funkcí.

Také část zabývající se dolováním z dat skýtá velké množství možností pro její pokračování. Jako jeden z příkladů lze uvést rozšíření založené na datech získaných prostřednictvím dokumentu *User Tune*. Konkrétně se jedná o internetové stránky, které by také mimo jiné nabízely prodej hudby. Jako vzor lze brát stránku *Last.fm* [27], která jednotlivým uživatelům vytváří hudební profily na základě poslouchaných skladeb. Navrhované rozšíření by získávalo data pro tvorbu profilu z informací, které by byly šířeny prostřednictvím sítě Jabber a dokumentu *User Tune*. Využití dalšího standardu *User Location* by nahradilo manuální nastavení geografické polohy uživatele za automatické a tím by přineslo možnosti pro nabízení cílené reklamy.

S dokumentem *User Location* souvisí i další navrhované pokračování. Konkrétně se jedná o možný návrh nového XEP standardu, který by využíval určení geografické polohy například k naplánování cesty. Rozšíření by umožňovalo konkrétním lokalizačním souřadnicím přiřazovat důležité body, informace a poznámky. Ty by byly ihned rozeslány způsobem, jakým pracují dosavadní rozšíření postavená na *Personal Eventing Protocol* [12]. Využití by našlo ať už v záchranných misích nebo jen k šíření zajímavostí mezi uživateli.

V měsíci březen tohoto roku vyšla nová verze základních dokumentů RFC. V této práci již ale z těchto standardů nejsou zahrnuty žádné změny. V době vydání se totiž práce dostávala do konečné fáze vývoje. Proto změny přinesené novými protokoly mohly být taktéž uplatněny při případném pokračování této práce.

## Kapitola 5

# Vyhodnocení výsledků

Pátá kapitola je zaměřena na vyhodnocení výsledků a prezentaci zjištěných informací. V první podkapitole jsou uvedeny poznatky zjištěné z ručního prozkoumávání získaných dat, která jsou uložena v databázi. Z velké části je zde popsán jak vznik nepřesností, tak i jejich případné následné odstranění. V závěru této podkapitoly jsou shrnuta kvantitativní fakta zaměřená na databázi jako je například počet záznamů a velikost samotné databáze.

Druhá podkapitola se zabývá samotným procesem dolování z dat, který je prováděn pomocí nástroje RapidMiner [28]. Je postavena na informacích získaných z druhé a třetí kapitoly. Prostřednictvím tabulky 5.2, jejíž celá podoba se nachází v příloze F, jsou prezentovány zjištěné informace. Jedná se o procentuální shodu uživatelů mezi sebou, která je založena na datech získaných z jejich statusů v Jabber síti.

### 5.1 Manuální rozbor dat

Data, která jsou získaná z Jabber komunikace, jsou uložena v databázi PostgreSQL. V průběhu celého vývoje této práce byl obsah databáze důkladně kontrolován a případné problémy řešeny v nejbližším možném termínu. I tak se ale v databázi nacházelo několik chyb, které při následné analýze dokázaly mírně znepřesnit výsledky. Jako příklad lze uvést zaznamenávání rozšířeného statusu o geografické poloze do tabulky geoloc. V případě, že uživatel současně používal pouze jedno připojení, bylo vše v pořádku. Menší kolize nastala v situaci, kdy bylo více současných připojení. Každé rozesílalo své informace, které byly stejné, a tím se v jeden čas v databázi objevila redundantní data. Tato nepřesnost byla po manuálním rozboru dat objevena a odstraněna pomocí časových razítek.

Největší problém při sbírání dat je tvořen díky aplikacím, které jsou využívány uživateli. Jejich neúplná implementace XEP dokumentů a vzájemná nepřesná spolupráce s jinými programy tvoří velké množství matoucích dat, které se taktéž nacházejí v databázi. Tento problém se v největší míře objevuje v tabulce tune, jejímž cílem je uchování informací o hudbě, kterou uživatelé poslouchají. Podle dokumentu User Tune [18], by po skončení přehrávání hudby měl klient rozeslat do „světa“ zprávu, jejíž obsah, nesoucí informace o právě přehrávané hudbě, bude prázdný. Tím je sděleno ukončení přehrávání hudby. V mnoha případech se ale takto neděje. Zprávu, která obsahuje informace o poslední přehrávané hudbě, uživatel rozešle jako poslední. Tudíž z analýzy dat obsažených v databázi je mylně předpokládáno, že klient hudbu stále poslouchá.

Samotný sběr informací započal již v minulém roce v měsíci prosinec. V této době však ještě nebyla implementována všechna rozšíření a formát i obsah vybraných tabulek



byl později modifikován. Základní informace o uživatelském stavu, které jsou uchovávány v tabulce presence, již ale výraznou změnou neprošly a tudíž jsou sbírány od samého počátku běhu aplikace. V obsahu této entity pouze nejsou zahrnuty vybrané dny, kdy byl klient mimo provoz.

Databáze uchovává informace pouze o počtu kolem 40 uživatelů. Tento fakt výrazně ovlivnil přesnost a rozsah nasbíraných dat. Taktéž malé množství podpory rozšířených statusů ze strany uživatelů a jejich aplikací značně ovlivnilo obsah databáze. Ať už do počtu záznamů nebo umístění informací do správných entit. S tímto taktéž úzce souvisí rozšíření User Tune, jehož nesprávné použití je popsáno v závěrečné části druhé kapitoly.

Za dobu, která přesahuje 4 měsíce, bylo nasbíráno téměř 200 MB dat. Například tabulka, která uchovává data pro data mining popsaný v následující podkapitole, obsahuje přes 115 tisíc záznamů. Na druhou stranu entity zaměřené na rozšířené statusy, konkrétně to jsou geoloc, tune, mood a activity, disponují s menším počtem záznamů, které se počítají pouze ve stovkách. Tento omezený počet, jak již bylo uvedeno, je výsledkem menší podpory těchto rozšíření ze strany klientů.

Tabulka 5.1 souhrnně prezentuje nasbíraná data v jednotlivých kategoriích rozšíření. Prvními zjištěnými údaji je tvořen druhý sloupec, kterým je shrnut počet uživatelů podporujících daný standard. Jak druhý tak i třetí sloupec se skládá ze dvou čísel, která jsou vzájemně oddělena znakem lomítko /. První údaj uvádí počet užitečných záznamů z celkového počtu. Konkrétně v druhém sloupci to je počet uživatelů, kteří šířili užitečné informace, k počtu uživatelů, jež šířili jakákoliv data, tedy i prázdná. Tabulku uzavírá přehled průměrného počtu záznamů na jednoho uživatele.

Rozšíření	Uživatelé/z	Záznamy/z	Průměr
User Tune	10 %(4)/44 %(18)	2258/4789	564
User Location	12 %(5)/12 %(5)	122/276	24
User Mood	20 %(8)/31 %(13)	550/1367	63
User Activity	10 %(4)/20 %(8)	404/1136	101

Tabulka 5.1: Přehled manuálního rozboru dat.

Manuální rozbor dat probíhal i při tvorbě webové prezentace, která se zaměřuje na praktické využití získaných dat. Lze ji nalézt na stránce <http://pcpanel-1205.fit.vutbr.cz/jabinfo/> na jejíž tvorbě se z velké části podílel vedoucí této práce Ing. Jozef Mlích. Obsah této ukázky slouží i jako názorný manuál, jak by se s daty mohlo dále nakládat.

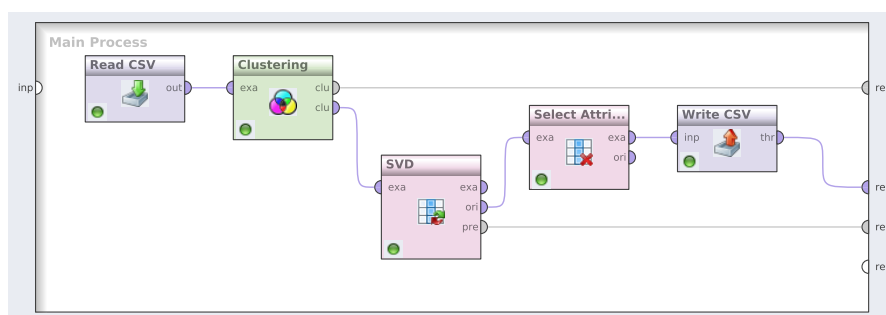
Následující podkapitola se zaměřuje na vyhodnocení údajů z tabulky presence. Toto rozhodnutí je postaveno na ručním prozkoumání získaných dat, které je shrnuto v tabulce 5.1. Jsou-li použita data z entit mood a activity je na první pohled patrné, že uživatelé využívající tyto služby, je ve své aplikaci nastavili pouze jedenkrát. Jelikož nastavení těchto rozšíření neprobíhá automaticky, nelze získané informace považovat za věrohodné. Pro samotnou analýzu by byl potřebný mnohonásobně větší objem dat, než jaký by byl nutný například pro tabulky tune a geoloc. Po prozkoumání tabulky geoloc bylo vyvozeno, že získané geografické polohy od klientů nejsou dále použitelné.

## 5.2 Programové vyhodnocení

V této části vyhodnocení jsou využita data získaná robotem z Jabber komunikace, která jsou uschována v databázi. Před samotným procesem dolování z dat, na řadu přichází

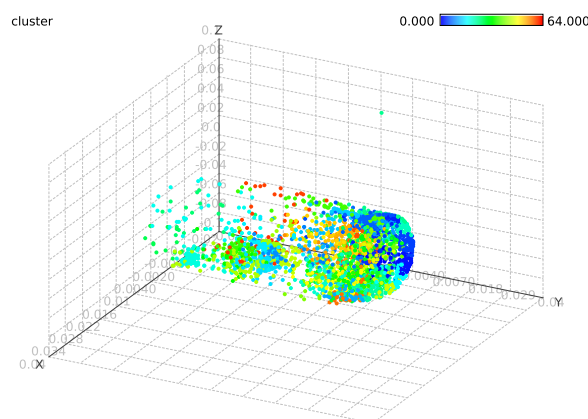
transformace popsaná v kapitole třetí, která data připraví. Je vytvořen formát, který je kompatibilní s programem RapidMiner [28]. Jak bylo uvedeno na konci předešlé kapitoly, rozšířené statusy pro proces dolování z dat nejsou vhodné. I z tohoto důvodu, ale hlavně díky velkému množství nasbíraných dat byla pro data mining využita tabulka presence, jejíž obsah je tvořen informacemi o změně statusů jednotlivých uživatelů.

Struktura samotného procesu vytvořeného v programu RapidMiner je zobrazena na obrázku 5.1. Jak je patrné je tvořena z několika bloků, kde každý plní určitou funkci. Vzájemné vztahy mezi jednotlivými oddíly procesu jsou reprezentovány úsečkami. První částí, která je nazvána *Read CSV*, jsou transformovaná data načtena z externího souboru do nástroje. Následující oddíl *Clustering* využívá algoritmus *k*-means, kterým se podrobně zabývá kapitola třetí. Počet shluků byl experimentováním nastaven na konečnou hodnotu 100. Blok nazvaný *SVD*<sup>1</sup> slouží ke snížení dimenzí jednotlivých vstupních vektorů. Pro



Obrázek 5.1: Struktura procesu v nástroji RapidMiner.

vizuální zobrazení bylo zvoleno snížení dimenzí z 288 až na hodnotu 3. Následujícím oddílem *Select Attribute*, který využívá data z bloku *Clustering*, jsou vybrány pouze dva atributy. Konkrétně to jsou jednotlivá JID uživatelů a jména shluků, do kterých byly zařazeny. Posledním blokem jsou tato data exportována do CSV souboru.



Obrázek 5.2: Výsledné shluky převedeny do 3D prostoru.

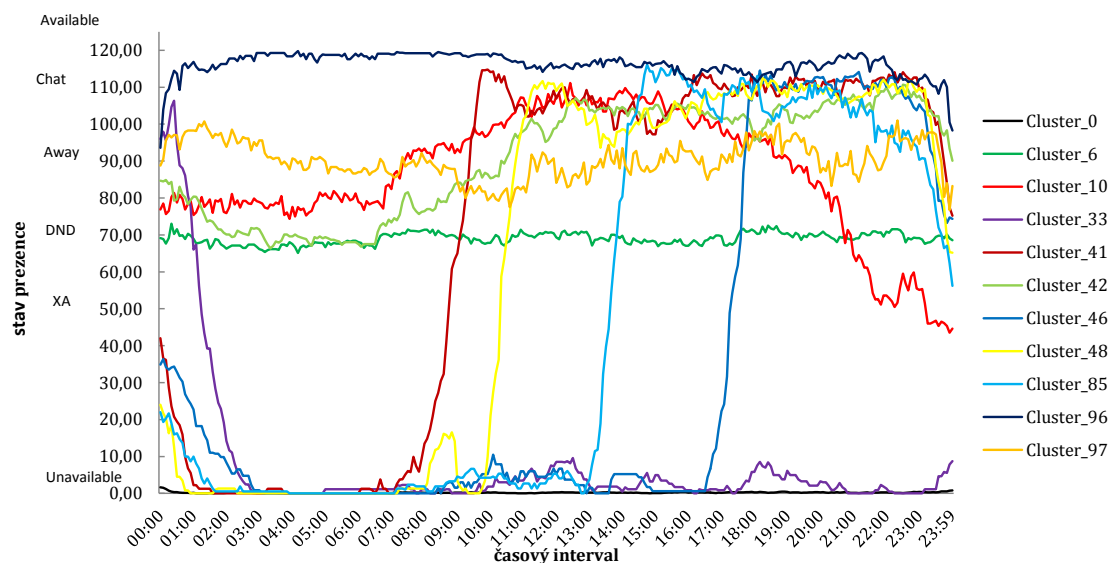
Obrázkem 5.2 je prezentováno rozmístění jednotlivých shluků v tří rozměrném prostoru.

<sup>1</sup>Singular Value Decomposition



Byl vygenerován automaticky blokem SVD zobrazeném na obrázku 5.1. Tento blok, jak již bylo uvedeno, redukuje  $n$ -rozměrný vstupní vektor do tří dimenzí. Snižování prostorového řádu bylo provedeno za účelem grafického znázornění jednotlivých shluků. Tyto výsledné množiny jsou od sebe navzájem odlišeny barvami.

Pomocí algoritmu  $k$ -means byli uživatelé rozčleněni podle svého chování do několika množin. Pouze jedenáct shluků obsahovalo více jak osmdesát záznamů a právě ty jsou prezentovány na obrázku 5.3. Je zde znázorněn graf znázorňující vztah mezi statusem a časem. Osa  $x$  je tvořena časovým intervalem v rozmezí jednoho dne a na ose  $y$  je vynesena aktuální stav neboli presence.



Obrázek 5.3: Průběh statusů shluků podle času.

Po provedení shlukování algoritmem  $k$ -means pomocí nástroje RapidMiner jsou uživatelé rozřazeni do shluků a uloženi do externího souboru. V této chvíli se celkové vyhodnocení dostává do závěrečné fáze. Na řadu přichází aplikace, pomocí které je vypočítáno zastoupení jednotlivých uživatelů ve shlucích.

Vzájemná procentuální shoda uživatelů mezi sebou je prezentována tabulkou 5.2. Kde hlavička, která je podbarvena šedě a umístěna na prvním řádku, znázorňuje vybrané uživatele. Jejich samotné JID je, za účelem úspory místa a zachování soukromí, nahrazeno čísly od jedné po dvaadvacet. Tato struktura je uvedena i v prvním sloupci. Je tedy vytvořena „matice“, jejíž hlavní diagonála není vyplněna. Obsah je tvořen procentuální shodou mezi jednotlivými uživateli. Buňky prezentující procentuální zastoupení větší než 80 % jsou barevně odlišeny. Zeleně zabarvená políčka znázorňují podobnost v intervalu od 80 do 89 procent a červená barva značí interval shody mezi 90 až 99.

Výše zmíněná tabulka 5.2 prezentuje pouze vybranou část uživatelů. Celý přehled je umístěn v příloze F v tabulce F.2. Jednotlivé nahrazení JID uživatelů za čísla je taktéž uvedeno v příloze F, konkrétně v tabulce F.1.

Údaje prezentované tabulkou 5.2 byly i manuálně zkontrolovány a to vedlo k závěru, že celý proces počínaje transformací až po samotné vyhodnocení je navržen správně. Je patrné, že ke 100 procentní shodě nedošlo ani u uživatelů, kteří využívají dva stejné účty. Tento důsledek je možné vysvětlit tím, že ne vždy byly oba účty ve stejném stavu a ve stejné chvíli

-	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
1	-	22	31	70	23	28	31	37	55	54	64	14	52	32	34	54	44	27	23	19	29	23
2	22	-	63	24	90	88	84	55	52	39	34	90	29	14	56	32	44	89	85	50	33	22
3	31	63	-	26	63	66	64	57	55	47	42	60	45	22	79	34	44	66	63	57	36	38
4	70	24	26	-	23	29	30	36	53	48	58	14	42	26	28	58	41	26	21	18	30	19
5	23	90	63	23	-	90	85	60	49	42	34	88	32	17	55	30	42	86	86	49	34	23
6	28	88	66	29	90	-	86	60	56	44	38	85	34	17	58	36	42	87	85	49	32	24
7	31	84	64	30	85	86	-	59	53	49	43	81	31	22	56	32	42	86	83	51	34	25
8	37	55	57	36	60	60	59	-	60	53	48	51	41	21	60	39	51	58	59	50	43	24
9	55	52	55	53	49	56	53	60	-	58	58	43	53	26	55	52	63	56	52	47	43	26
10	54	39	47	48	42	44	49	53	58	-	60	33	49	31	48	55	56	44	42	39	44	29
11	64	34	42	58	34	38	43	48	58	60	-	26	48	35	44	55	46	38	34	31	44	29
12	14	90	60	14	88	85	81	51	43	33	26	-	24	9	51	25	36	82	82	49	29	21
13	52	29	45	42	32	34	31	41	53	49	48	24	-	29	45	56	53	37	29	29	40	30
14	32	14	22	26	17	17	22	21	26	31	35	9	29	-	24	25	22	17	14	16	23	19
15	34	56	79	28	55	58	56	60	55	48	44	51	45	24	-	37	49	59	56	59	39	44
16	54	32	34	58	30	36	32	39	52	55	55	25	56	25	37	-	49	34	30	27	36	25
17	44	44	44	41	42	42	42	51	63	56	46	36	53	22	49	49	-	49	46	38	49	25
18	27	89	66	26	86	87	86	58	56	44	38	82	37	17	59	34	49	-	86	51	36	23
19	23	85	63	21	86	85	83	59	52	42	34	82	29	14	56	30	46	86	-	51	36	25
20	19	50	57	18	49	49	51	50	47	39	31	49	29	16	59	27	38	51	51	-	34	35
21	29	33	36	30	34	32	34	43	43	44	44	29	40	23	39	36	49	36	36	34	-	26
22	23	22	38	19	23	24	25	24	26	29	29	21	30	19	44	25	25	23	25	35	26	-

Tabulka 5.2: Procentuální shoda jednotlivých uživatelů.

zapínány. Taktéž se dospělo ke zjištění, které dokazuje, že ke správnému vyhodnocení by bylo potřeba, aby všichni uživatelé měli přibližně stejný počet záznamů obsahujících změnu statusu. K velké procentuální shodě například došlo u dvou klientů, kteří měli v databázi pouze tři záznamy, což nemá žádnou vypovídající hodnotu o jejich podobnosti.

## Kapitola 6

# Závěr

Cílem této bakalářské práce bylo vytvoření klienta pro síť Jabber pracující v reálném čase. Tato aplikace pracuje jako robot, který je schopen reagovat na vhodně zvolené příkazy a sbírat data z komunikace. Následujícím krokem bylo analyzování nashromážděných dat, se kterými byl dále proveden proces dolování. Výsledky jsou prezentovány pomocí tabulky, jejíž obsah je tvořen procentuální shodou jednotlivých uživatelů mezi sebou.

Po důkladném vypracování této práce je možné prohlásit, že všechny cíle kladené na tuto práci se podařilo splnit. Každému bodu zadání je zde věnován oddíl, ve kterém je daná problematika podrobně rozebrána. Prvnímu bodu zadání odpovídá celá druhá kapitola, která je tvořena popisem protokolu XMPP, na kterém je postavena komunikační služba Jabber. Také je zde rozebrána architektura sítě a základní stavební kameny XMPP protokolu, které jsou využívány sítí Jabber jako nástroj ke zprostředkování jednotlivých služeb. Druhý bod zadání úzce souvisí s tvorbou robota, který byl vytvořen za účelem získávání dat z Jabber komunikace. Zdrojové soubory této aplikace jsou přiloženy na kompaktním disku. Implementační částí se věnuje kapitola čtvrtá, kde je rozebrána jak celková programová architektura této práce tak její jednotlivé bloky. Prvním oddílem kapitoly páté je splnění třetí bod zadání. Jsou zde uvedeny poznatky, které byly získány z manuálního průzkumu dat. Z ručního prohlížení dat byly vyvozeny důsledky, které tvořily základ pro splnění bodu čtvrtého. Tomuto a následnému bodu zadání je věnována zbývající část páté kapitoly, která je doplněna o přílohy nacházející se na konci tohoto dokumentu. Pro automatické zkoumání dat byl využit nástroj RapidMiner, který je společně s dalšími vybranými programy popsán na konci kapitoly třetí. Její obsah je zaměřen na popis procesu dolování z dat a na jeho začlenění do celého bloku získávání znalostí z databází. Poslednímu bodu zadání je vyhrazena část v kapitole čtvrté, která se zabývá případným pokračováním této práce.

Jak je uvedeno výše veškeré body zadání byly splněny. Jako rozšíření zadání bylo zpřístupnění možnosti využít jak robota tak i data uložená v databázi pro prezentování vlastního statusu na webových stránkách. Příklad tohoto použití je přiložený na kompaktním disku a také je k dispozici k nahlédnutí na internetové stránce <http://pcpanel-1205.fit.vutbr.cz/jabinfo/>.

Pro nedostatečný počet uživatelů, kteří doposud využívali tohoto robota, a zároveň by měli správně nastavené „rozesílání“ rozšířených statusů vedlo k následujícímu rozhodnutí. Větší důraz byl kladen na vývoj samotného robota a sbírání i ukládání dat do databáze, než na následný proces dolování z dat. Za období čtyř měsíců bylo celkem nasbíráno kolem 200 MB dat a přes 100 tisíc záznamů.

# Literatura

- [1] Adams, D.: *Programming jabber*. Sebastopol: O'Reilly, první vydání, 2002, 455 s., iISBN 05-960-0202-5.
- [2] Berka, P.: *Dobývání znalostí z databází*. Praha: Academia, první vydání, 2003, 366 s., iISBN 80-200-1062-9.
- [3] Bramer, M.: *Principles of Data mining*. London: Springer, první vydání, 2007, 343 s., iISBN 18-462-8765-0.
- [4] Fayyad, U. M.; Smyth, P.: *Advances in knowledge discovery and data mining*. California: MIT Press, první vydání, 1996, 611 s., iISBN 02-625-6097-6.
- [5] Han, J.; Kamber, M.: *Data mining : concepts and techniques*. San Francisco: Morgan Kaufmann Publisher, druhé vydání, 2006, 770 s., iISBN 15-586-0901-6.
- [6] Hildebrand, J.; Saint-Andre, P.: XEP-0080: User Location. [online], 15-09-2009, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0080.html>
- [7] Hildebrand, J.; Saint-Andre, P.; Tronçon, R.; aj.: XEP-0115: Entity Capabilities. [online], 26-02-2008, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0115.html>
- [8] Hruška, T.: *Informační systémy : IIS/PIS*. Brno: Fakulta informačních technologií, 2008, 14733 s.
- [9] Kolektiv autorů: Extensible Markup Language (XML) 1.0. [online], 26-11-2008, [cit. 16. května 2011].  
URL <http://www.w3.org/TR/2008/REC-xml-20081126/>
- [10] Kosek, J.: *XML pro každého : podrobný průvodce*. Praha: Grada, první vydání, 2000, 163 s., iISBN 80-716-9860-1.
- [11] Meijer, R.; Saint-Andre, P.: XEP-0108: User Activity. [online], 29-10-2008, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0108.html>
- [12] Millard, P.; Saint-Andre, P.; Meijer, R.: XEP-0060: Publish-Subscribe. [online], 12-07-2010, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0060.html>

- [13] Mizzi, S.; Saint-Andre, P.: XEP-0292: vCard4 Over XMPP. [online], 02-26-2008, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0292.html>
- [14] Moore, D.; Wright, W.: *Jabber developer's handbook*. Indianapolis: Sams Publishing, první vydání, 2004, 487 s., iISBN 06-723-2536-5.
- [15] Nemrava, M.; Pospíšil, J.: Dolování dat a jeho aplikace. [online], 2006, [cit. 16. května 2011].  
URL [http://www.spatial.cs.umn.edu/paper\\_ps/dmchap.pdf](http://www.spatial.cs.umn.edu/paper_ps/dmchap.pdf)
- [16] Saint-Andre, P.: XEP-0054: vcard-temp. [online], 07-16-2008, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0054.html>
- [17] Saint-Andre, P.: XEP-0092: Software Version. [online], 02-15-2007, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0092.html>
- [18] Saint-Andre, P.: XEP-0118: User Tune. [online], 30-01-2008, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0118.html>
- [19] Saint-Andre, P.: XEP-0153: vCard-Based Avatars. [online], 16-08-2006, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0153.html>
- [20] Saint-Andre, P.: XEP-0199: XMPP Ping. [online], 03-06-2009, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0199.html>
- [21] Saint-Andre, P.; Meijer, R.: XEP-0107: User Mood. [online], 29-10-2008, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0107.html>
- [22] Saint-Andre, P.; Smith, K.: XEP-0163: Personal Eventing Protocol. [online], 12-07-2010, [cit. 16. května 2011].  
URL <http://xmpp.org/extensions/xep-0163.html>
- [23] Saint-André, P.: Extensible Messaging and Presence Protocol (XMPP): Core. [online], 10-2004, [cit. 16. května 2011].  
URL <http://tools.ietf.org/html/rfc3920>
- [24] Saint-André, P.: Extensible Messaging and Presence Protocol (XMPP): Instant Messaging and Presence. [online], 10-2004, [cit. 16. května 2011].  
URL <http://tools.ietf.org/html/rfc3921>
- [25] Saint-André, P.; Smith, K.; Troncon, R.: *XMPP : the definitive guide : building real-time applications with jabber technologies*. Sebastopol: O'Reilly, první vydání, 2009, 287 s., iISBN 978-059-6521-264.
- [26] WWW Stránky: Database Systems. [online], 2007, [cit. 16. května 2011].  
URL [http://www.cs.uct.ac.za/mit\\_notes\\_devel/Database/Lates%t/index.html](http://www.cs.uct.ac.za/mit_notes_devel/Database/Lates%t/index.html)

- [27] WWW Stránky: Last.fm. [online], 02-11-2009 [cit. 16. května 2011].  
URL <http://www.last.fm/>
- [28] WWW Stránky: RapidMiner. [online], 2011, [cit. 16. května 2011].  
URL <http://rapid-i.com/content/view/181/190/>
- [29] Řezánková, H.; Húsek, D.; Snášel, V.: *Shluková analýza dat*. Praha: Professional Publishing, druhé vydání, 2009, 218 s., iISBN 978-808-6946-818.

# Příloha A

## Obsah CD

<code>./doc/</code>	- programová dokumentace
<code>./src/datamining</code>	- zdrojové soubory pro data mining
<code>./src/jabinfo</code>	- zdrojové soubory robota
<code>./src/transformation</code>	- zdrojové soubory transformace
<code>./tex/</code>	- zdrojová data technické zprávy
<code>./www/</code>	- zdrojová data webové prezentace
<code>./bp-xsendl00.pdf</code>	- technická zpráva
<code>./MANUAL</code>	- manuál k programu jabinfo
<code>./README</code>	- základní informace

## Příloha B

# Slovník zkratek

**Base64** — převedení binárních dat pomocí znaků ASCII.

**CSV** (Comma-separated values) — souborový formát určený pro výměnu dat z tabulek.

**GPG** (GNU Privacy Guard) — slouží k podepisování a kontrolu podpisu různých dat.

**ID3v1** — datový formát obsahující informace o skladbě.

**IM služby** (Instant messaging) — umožňují posílat zprávy a komunikovat mezi uživateli.

**Jabber** — viz XMPP.

**JEP** (Jabber Enhancement Proposal) — starší název pro XEP.

**JID** (Jabber ID) — jednoznačný identifikátor v síti Jabber.

**SASL** (Simple Authentication and Security Layer) — metoda sloužící ověřování klientů na serverech.

**Stanza** — XML stream.

**SQL** (Structured Query Language) — jazyk používaný pro komunikaci s databázemi.

**TLS** (Transport Layer Security) — kryptografický protokol, poskytuje zabezpečenou komunikaci na internetu.

**TSQL** (Transact-SQL) — jazyk pro temporální databáze.

**vCard** — elektronická osobní vizitka.

**XEP** (XMPP Extension Protocol) — rozšiřuje protokol RFC.

**XML** (Extensible Markup Language) — univerzální značkový jazyk.

**XMPP** (Extensible Messaging and Presence Protocol) — protokol pro doručování zpráv.



## Příloha C

# Stanza - základní schéma

Přehled základních elementů, které jsou využívány při Jabber komunikaci. Struktura jednotlivých částí stanzy ukazuje pouze prvky relativní k této práci. Pomocí hranatých závorek je znázorněna množina, ze které musí být vybrán právě jeden prvek. Na místě uvozovek se očekává jakákoliv povolená hodnota.

### Iq

---

```
1      <iq from=""
2          to=""
3          type="[ get , set , result , error ]"
4          id=""
5          Namespace
6      </iq>
```

---

Algoritmus C.1: Popis elementu *iq*.

### Message

---

```
1      <message from=""
2          to=""
3          type="[ normal , chat , groupchat , headline , error ]"
4          id=""
5          <body> </body>
6          <x xmlns="jabber:x:event">
7              [ Offline , Delivered , Displayed , Composing ]
8          <subject> </subject>
9          <thread> </thread>
10         <error> </error>
11         <x> </x>
12     </message>
```

---

Algoritmus C.2: Popis elementu *message*.

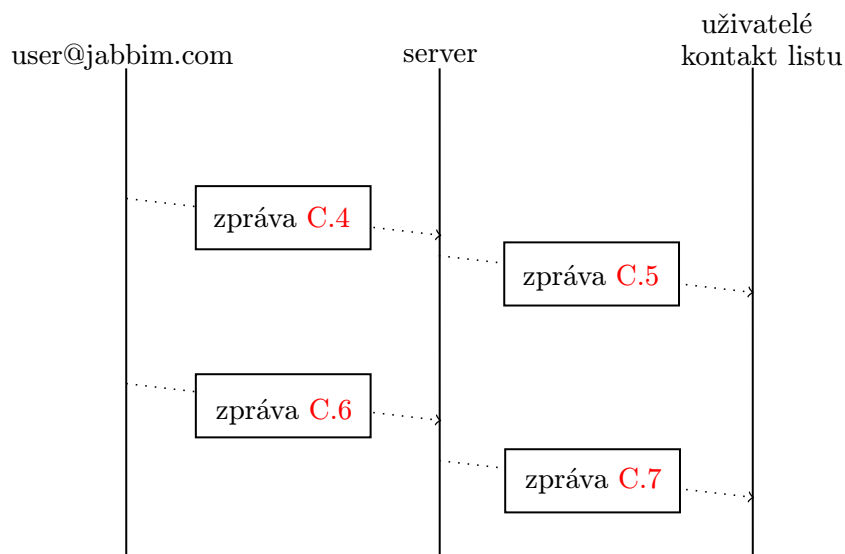
## Presence

```
1  <presence from=""  
2      to=""  
3      type="[available , unavailable , probe , subscribe ,  
4            unsubscribe , subscribed , unsubscribed , error]"  
5      id=""  
6  <show>  
7      [away , chat , dnd , normal , xa]  
8  </show>  
9  <status>    </status>  
10 <priority>   </priority>  
11 <error>      </error>  
12 </presence>
```

Algoritmus C.3: Popis elementu *presence*.

## Přehled průběhu rozšíření

Ukázka celého příkladu šíření statusu pomocí rozšíření *User Tune*. Uživatel *user* poslouchá hudbu a informuje server zasláním zprávy zobrazené v příkladu C.4.



Obrázek C.1: Ukázka „šíření“ *User Tune*.

---

```

1  <iq from="user@jabber.com" type="set" id="pub1">
2    <pubsub xmlns="http://jabber.org/protocol/pubsub">
3      <publish node="http://jabber.org/protocol/tune">
4        <item>
5          <tune xmlns="http://jabber.org/protocol/tune">
6            <artist>Daniel Landa</artist>
7            <length>255</length>
8            <source>Nigredo</source>
9            <title>1968</title>
10           <track>5</track>
11         </tune>
12       </item>
13     </publish>
14   </pubsub>
15 </iq>

```

---

Algoritmus C.4: Informování serveru o právě přehrávané hudbě.

Server obdrží informace od klienta *user* prostřednictvím zprávy, jejíž obsah je naplněn daty spojenými s přehrávanou hudbou. Pomocí elementu *message* ji přepošle všem uživatelům ze seznamu kontaktů uživatele *user*, kteří se pro odběr těchto typů zpráv zaregistrovali. Struktura této zprávy je prezentována na příkladu C.5.

---

```

1  <message from="user@jabber.com" type="set"
2    to="jabinfo@jabber.com/bot" id="pub1">
3    <event xmlns="http://jabber.org/protocol/pubsub#event">
4      <items node="http://jabber.org/protocol/tune">
5        <item>
6          <tune xmlns="http://jabber.org/protocol/tune">
7            <artist>Daniel Landa</artist>
8            <length>255</length>
9            <source>Nigredo</source>
10           <title>1968</title>
11           <track>5</track>
12         </tune>
13       </item>
14     </items>
15   </event>
16 </message>

```

---

Algoritmus C.5: Server informuje uživatele podporující rozšíření o stavu *user@jabber.com*.

Zpráva o přehrávané hudbě je také přeposílána všem otevřeným spojením uživatele *user*, ukázáno na příkladu C.6.

---

```
1  <message from="user@jabbbim.com" type="set"
2    to="user@jabbbim.com/дома" id="pub2">
3    <event xmlns="http://jabber.org/protocol/pubsub#event">
4      <items node="http://jabber.org/protocol/tune">
5        <item>
6          <tune xmlns="http://jabber.org/protocol/tune">
7            <artist>Daniel Landa</artist>
8            <length>255</length>
9            <source>Nigredo</source>
10           <title>1968</title>
11           <track>5</track>
12         </tune>
13       </item>
14     </items>
15   </event>
16 </message>
```

---

Algoritmus C.6: Server přepośle informace o přehrávané hudbě všem otevřeným spojením uživatele *user@jabbbim.com*.

Přestane-li uživatel *user* poslouchat/vysílat informace o přehrávané hudbě, provede to pomocí zprávy ukázané na příkladu C.7. Zpráva typu *iq*, ve které je položka *tune*, nesoucí informace o skladbě, je prázdná.

---

```
1  <iq from="user@jabbbim.com/prace" type="set" id="pub1">
2    <pubsub xmlns="http://jabber.org/protocol/pubsub">
3      <publish node="http://jabber.org/protocol/tune">
4        <item>
5          <tune xmlns="http://jabber.org/protocol/tune"/>
6        </item>
7      </publish>
8    </pubsub>
9  </iq>
```

---

Algoritmus C.7: Uživatel ukončil „vysílání“ rozšířených zpráv o svém stavu.

Server informuje všechny účastníky odběru zprávou, která má položku *tune* prázdnou. Tak jak to prezentuje příklad C.8.

---

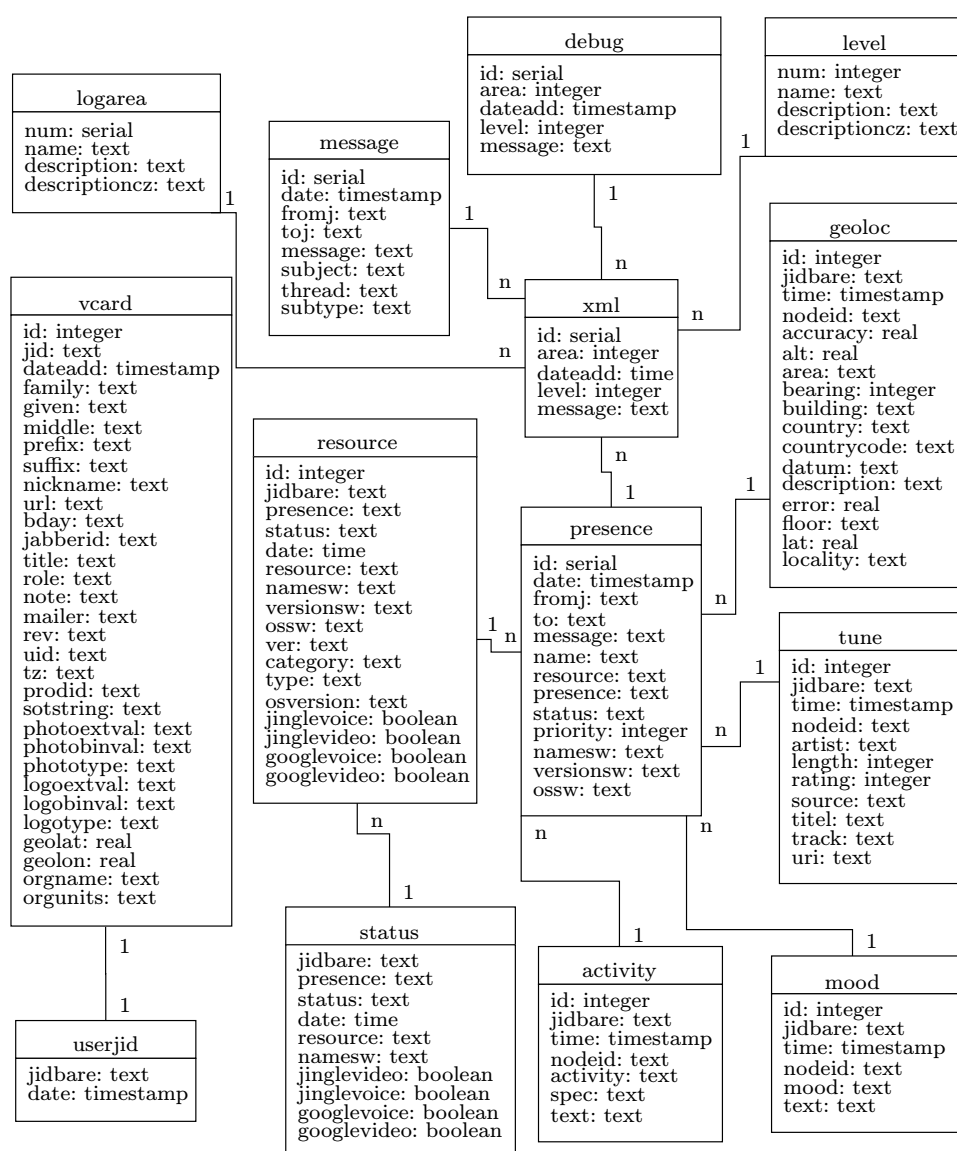
```
1  <message from="user@jabbbim.com"
2    to="jabinfo@jabbbim.com/bot">
3    <event xmlns="http://jabber.org/protocol/pubsub#event">
4      <items node="http://jabber.org/protocol/tune">
5        <item>
6          <tune xmlns="http://jabber.org/protocol/tune"/>
7        </item>
8      </items>
9    </event>
10 </message>
```

---

Algoritmus C.8: Server informuje klienty o ukončení šíření rozšířeného statusu.

## Příloha D

# Návrh databáze



Obrázek D.1: Struktura databáze

## Příloha E

# Přehled příkazů robota

Přehled základních zpráv obsahující příkazy pro robota a jejich vysvětlení je prezentováno v tabulce E.1. Kde první sloupec značí příkaz, který rozlišuje velká a malá písmena, a druhý jej krátce a výstižně popisuje. V případě, že uživatel nemá žádný údaj v databázi potřebný k vyhodnocení dotazu, je zaslána pouze informativní zpráva, které vzniklý problém vysvětluje. V případě zaslání neznámého dotazu, je uživateli automaticky odeslána zpráva s obsahem příkazu HELP.

příkaz	vysvětlení
HELP	Přehled příkazů a jejich popis.
INFO	Popis aplikace.
HISTORY	Výpis posledních 10 změn statusů.
TUNE	Výpis o přehrávané hudbě.
GEOLOC	Výpis o geografické poloze.
MOOD	Výpis o aktuální náladě.
ACTIVITY	Výpis o aktuální činnosti.
VCARD	Základní obsah VCard

Tabulka E.1: Přehled příkazů pro robota.

## Příloha F

### Přehled vztahů uživatelů

Tabulkou [F.1](#) je prezentován vztah mezi uživatelským JID a náhodně přiřazeným číslem, které jej v této práci zastupuje. Klientské jednoznačné identifikační názvy jsou zde zobrazeny bez části, která následuje za znakem zavináč, tedy bez názvu serveru. K tomuto kroku bylo přistoupeno ke snaze zajistit alespoň minimální stupeň anonymity.

JID	číslo	JID	číslo
portilo	1	mewoack	22
xsendl00	2	yac	23
hrabos	3	komi	24
kubira	4	drake127	25
mautinek	5	fletcher	26
danulka	6	boky	27
matmatmat	7	mtmccox	28
ybyzak	8	tux.martin	29
vecerniceverca	9	martin	30
pojir	10	tux.martin	31
vtheman	11	xapoh	32
jurij.h	12	matasx	33
imlich	13	javier.brat	34
joejoe	14	sychra	35
semal	15	ondrg	36
mr.aston	16	moen	37
mrazek.v	17	myler	38
theslayer	18	martix	39
kuba86	19	rezza	40
minimax	20	ales.lanik	41
miklik	21		

Tabulka F.1: Zobrazení uživatelského JID na čísla.

Tabulka [F.2](#), která se nachází níže, plně zobrazuje procentuální shodu mezi jednotlivými uživateli. Čísla v prvním řádku a sloupci odpovídají přiřazení prezentovanému v tabulce [F.1](#). Barevně odlišené buňky tabulky vyzdvihují větší procentuální shodu než jaká je u ostatních hodnot. Barvě zelené odpovídá interval od 80 % – 89 % a pro barvu červenou 90 % – 99 %.

-	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41
1	-	22	31	70	23	28	31	37	55	54	64	14	52	32	34	54	44	27	23	19	29	23	60	52	51	25	29	53	33	34	35	61	49	40	53	32	23	27	14	26	22
2	22	-	63	24	90	88	84	55	52	39	34	90	29	14	56	32	44	89	85	50	33	22	31	36	43	36	29	34	53	61	60	33	29	25	49	81	84	48	91	77	92
3	31	63	-	26	63	66	64	57	55	47	42	60	45	22	79	34	44	66	63	57	36	38	36	45	48	53	51	37	51	58	58	37	36	29	55	69	65	57	60	68	63
4	70	24	26	-	23	29	30	36	53	48	58	14	42	26	28	58	41	22	18	30	19	45	42	42	21	23	51	33	36	55	42	36	47	29	23	15	25	21			
5	23	90	63	23	-	90	85	60	49	42	34	88	32	17	55	30	42	86	86	49	34	23	29	36	41	34	31	35	53	58	58	34	29	27	48	75	82	49	88	79	91
6	28	88	66	29	90	-	86	60	56	44	38	85	34	17	58	36	42	87	85	49	32	24	34	37	42	34	31	34	51	57	56	38	34	31	51	77	84	49	85	90	86
7	31	84	64	30	85	86	-	59	53	49	43	81	31	22	56	32	42	86	83	51	34	25	33	44	49	36	32	39	53	62	61	41	37	36	54	79	82	51	81	80	83
8	37	55	57	36	60	60	59	-	60	53	48	51	41	21	60	39	51	58	59	50	43	24	41	42	51	34	33	42	60	66	66	50	42	31	55	60	55	56	51	58	54
9	55	52	55	53	49	56	53	60	-	58	58	43	53	26	55	52	63	56	52	47	43	26	60	51	60	40	37	56	62	62	51	43	36	65	56	51	52	44	53	49	
10	54	39	47	48	42	44	49	53	58	-	60	33	49	31	48	55	56	44	42	39	44	29	58	61	62	42	45	61	44	46	47	59	54	47	66	50	45	34	48	41	
11	64	34	42	58	34	38	43	48	58	60	-	26	48	35	44	55	46	38	34	31	44	29	58	61	60	36	41	59	41	44	44	64	55	48	58	44	32	38	27	39	34
12	14	90	60	14	88	85	81	51	43	33	26	-	24	9	51	25	36	82	82	49	29	21	24	29	34	32	27	26	47	54	53	27	24	23	42	73	81	46	99	75	88
13	52	29	45	42	32	34	31	41	53	49	48	24	-	29	45	56	53	37	29	29	40	30	66	50	50	34	38	58	34	37	38	43	43	26	49	38	31	24	42	32	
14	32	14	22	26	17	17	22	21	26	31	35	9	29	-	24	25	22	17	14	16	23	19	26	42	32	19	25	27	15	16	16	24	25	26	25	20	16	10	21	15	
15	34	56	79	28	55	58	56	60	55	48	44	51	45	24	-	37	49	59	56	59	39	44	38	50	55	59	38	51	57	57	34	42	34	56	62	56	60	51	58	56	
16	54	32	34	58	30	36	32	39	52	55	55	25	56	25	37	-	49	34	30	27	36	25	58	47	48	32	33	58	33	38	53	48	38	56	41	34	29	25	38	32	
17	44	44	44	41	42	42	42	51	63	56	46	36	53	22	49	49	-	49	46	38	49	25	49	51	60	36	34	58	59	58	58	40	34	28	54	52	46	44	36	49	42
18	27	89	66	26	86	87	86	58	56	44	38	82	37	17	59	34	49	-	86	51	36	23	36	43	49	36	31	38	55	64	64	36	31	29	52	82	84	49	83	81	88
19	23	85	63	21	86	85	83	59	52	42	34	82	29	14	56	30	46	86	-	51	36	25	34	36	45	36	33	34	57	65	65	36	32	27	45	77	83	56	82	77	86
20	19	50	57	18	49	49	51	50	47	39	31	49	29	16	59	27	38	51	51	-	34	35	27	40	44	52	42	27	49	51	51	30	36	36	44	52	51	51	49	51	52
21	29	33	36	30	34	32	34	43	43	44	44	29	40	23	39	36	49	36	36	34	-	26	37	45	51	36	36	42	47	43	43	36	31	33	37	36	36	29	35	36	
22	23	22	38	19	23	24	25	24	26	29	21	30	19	44	25	25	23	25	35	26	-	26	30	31	70	57	25	23	23	23	25	44	36	27	26	24	35	21	26	25	
23	60	31	36	45	29	34	33	41	60	58	58	24	66	26	38	58	49	36	34	27	37	26	-	56	55	31	36	56	34	38	40	59	52	42	56	40	37	35	25	40	32
24	52	36	45	42	36	37	44	42	51	61	61	29	50	42	50	47	51	43	36	40	45	30	56	-	69	45	44	49	43	47	48	51	47	47	55	47	34	37	30	44	36
25	51	43	48	42	41	42	49	51	60	62	60	34	50	32	55	48	60	49	45	44	51	31	55	69	-	48	44	51	53	58	59	55	48	46	55	53	42	48	34	44	44
26	25	36	53	21	34	34	36	34	40	42	36	32	34	19	59	32	36	36	36	52	36	70	31	45	48	-	66	31	37	38	38	34	50	40	39	39	35	45	33	36	38
27	29	29	51	23	31	31	32	33	37	45	41	27	38	25	59	33	34	31	33	42	36	57	36	44	44	66	-	38	31	31	31	36	51	42	38	36	34	42	27	33	31
28	53	34	37	51	35	34	39	42	56	61	59	26	58	27	38	58	38	34	27	42	25	56	49	51	31	38	-	43	40	40	54	49	43	59	40	43	35	27	42	33	
29	33	53	51	33	53	51	53	60	62	44	41	47	34	15	51	33	59	55	57	49	47	23	34	43	53	37	31	43	-	80	79	43	31	25	49	55	48	55	47	49	50
30	34	61	58	36	58	57	62	66	62	46	44	54	37	16	57	38	58	64	65	51	43	23	38	47	58	38	31	40	80	-	98	45	36	29	52	65	56	54	55	58	59
31	35	60	58	36	58	56	61	66	62	47	44	53	38	16	57	38	58	64	65	51	43	23	40	48	59	38	31	40	79	98	-	45	36	30	52	65	55	54	54	58	58
32	61	33	37	55	34	38	41	50	51	59	64	27	43	24	34	53	40	36	36	30	36	25	59	51	55	34	36	54	43	45	45	-	65	51	63	41	32	37	27	35	33
33	49	29	36	42	29	34	37	42	43	54	55	24	43	25	42	48	34	31	32	36	31	44	52	47	48	50	51	49	31	36	65	-	56	53	36	29	40	25	31	31	
34	40	25	29	36	27	31	36	31	36	47	48	23	36	26	34	38	28	29	27	36	31	36	42	47	46	40	42	43	25	29	30	51	56	-	49	29	27	33	23	29	29
35	53	49	55	47	48	51	54	55	65	66	58	42	49	25	56	56	54	52	45	44	33	27	56	55	55	39	38	59	49	52	52	63	53	49	-	58	49	51	43	58	47
36	32	81	69	29	75	77	79	60	56	50	44	73	38	20	62	41	52	82	77	52	37	26	40	47	53	39	36	40	55	65	41	36	29	58	-	77	51	73	81	78	
37	23	84	65	23	82	84	82	55	51	45	32	81	38	16	56	34	46	84	83	51	36	24	37	34	42	35	34	43	48	56	55	32	29	27	49	77	-	47	81	82	84
38	27	48	57	23	49	49	51	56	52	45	38	46	31	16	60	29	44	49	56	51	36	35	37	48	45	42	35	55	54	54	37	40	33	51	51	47	-	46	49	49	
39	14	91	60	15	88	85	81	51	44	34	27	99	24	10	51	25	36	83	82	49	29	21	25	30	34	33	27	27	47	55	54	27	25	23	43	73	81	46	-	75	89
40	26	77	68	25	79	80	80	58	53	48	39	75	42	21	58	38	49	81	77	51	35	26	44	44	36	33	42	49	58	58	35	31	29	58	81	82	49	75	-	79	
41	22	92	63	21	91	86	83	54	49	41	34	88	32	15	56	32	42	88	86	52	36	25	44	38	31	32	36	44	38	31	59	58	33	31	29	47	78	84	49	89	79