

VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

BRNO UNIVERSITY OF TECHNOLOGY

FAKULTA INFORMAČNÍCH TECHNOLOGIÍ
ÚSTAV POČÍTAČOVÝCH SYSTÉMŮ

FACULTY OF INFORMATION TECHNOLOGY
DEPARTMENT OF COMPUTER SYSTEMS

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 POČÍTAČOVÝCH SYSTÉMŮ

FACULTY OF INFORMATION TECHNOLOGY
DEPARTMENT OF COMPUTER SYSTEMS

DATAMINING Z JABBERU

DATAMINING FROM JABBER

BAKALÁŘSKÁ PRÁCE

BACHELOR'S THESIS

AUTOR PRÁCE

AUTHOR

JAROSLAV SENDLER

VEDOUCÍ PRÁCE

SUPERVISOR

Ing. JOZEF MLÍCH

BRNO 2011

Abstrakt

Předmětem této bakalářské práce bylo seznámení se s problematikou komunikace přes Jabber síť, která zde byla rozebrána. Konkrétním cílem bylo vytvoření jednoduchého Jabber klienta, který by byl schopen získávat statistická data. Nashromážděná data sloužila pro pozdější analýzu a grafickou reprezentaci informací z nich získaných.

Abstract

The objective of this thesis was acquaint oneself with problems of communication via Jabber network, which was also analyzed. The specific objective was to create a simple Jabber's client which would be able to obtain statistical data. The collected data was used for analysis and graphic representation of information.

Klíčová slova

Jabber, XMPP, robot, datamining, dolování dat, RapidMiner.

Keywords

Jabber, XMPP, robot, datamining, 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

6. 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

1 Úvod	2
2 XMPP	4
2.1 Architektura	4
2.2 XML	7
2.3 Stanza	8
2.4 Rozšíření	10
3 Data mining	13
3.1 Transformace dat	16
3.2 Metody dolování dat	19
3.3 Shlukování	20
3.4 Programy	23
4 Implementace	25
4.1 Architektura	25
4.2 Databáze	26
4.3 Robot	28
5 Pokračování práce	30
6 Vyhodnocení výsledků	31
6.1 Manuální rozbor dat	31
7 Závěr	32
A Obsah CD	35
B Manual	36
C Konfigurační soubor	37
D Slovník zkratk	38
E Stanza - základní schéma	40
F Návrh databáze	45
G Přehled klientů a jejich rozšíření	46

1 Kapitola 1

2 Úvod

3 V současné době k nejrozšířenějším psaným dorozumívacím prostředkům jsou řazeny real-
4 time komunikační sítě. Samozřejmě tento jev je k vidění až konce minulého desetiletí. Psaná
5 komunikace je již několik století využívána jako prostředek k dorozumívání se mezi lidmi. Za
6 toto období prošla výrazným pokrokovým vývojem, proto u ní lze nalézt mnoho časových
7 milníků, které ji výrazně ovlivnily. Jako příklad jednoho z prvních komunikačních kanálů lze
8 uvést běžce, kteří často nesli zprávy na vzdálenosti několika kilometrů. Dalším výrazným
9 prvkem ve vývoji dorozumívacích prostředků bylo zavedení pošty a objevení telegrafu.

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

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

26 Předmětem této bakalářské práce bylo seznámení se s problematikou komunikace probí-
27 hající přes Jabber síť. Konkrétním cílem bylo vytvoření jednoduchého Jabber klienta, který
28 by byl schopen získávat statistická data. Nashromážděná data by dále sloužila pro pozdější
29 analýzu a grafickou reprezentaci informací z nich získaných. Tedy cílem této práce je získat
30 neznámé informace z real-time komunikační sítě Jabber.

31 Bakalářská práce je rozdělena do šesti kapitol. Kapitola **druhá** je tvořena popisem
32 protokolu XMPP, na kterém je postavena real-time komunikační služba Jabber. Je zde
33 popsána architektura sítě a základní stavební kameny XMPP protokolu, které jsou využí-
34 vány sítí Jabber jako nástroj k zprostředkování jednotlivých služeb. V závěru tohoto oddílu
35 je věnována část vybraným standardům, které rozšiřují základní XMPP protokoly. Jsou
36 zde uvedena pouze ta rozšíření, která byla po konzultaci s vedoucím práce, označena za
37 relevantní k této práci.

38 Obsah **třetí** kapitoly je zaměřen na popis procesu data mining. Zabývá se začleněním
39 této metody do komplexnějšího procesu, který je nazýván získávání znalostí z databází.
40 Další součást této kapitoly se zabývá nezbytnými kroky, při kterých jsou data, získaná z
41 Jabber komunikace, transformována. Při tomto procesu jsou data převedena do vhodné
42 podoby pro data minig. Následuje část, kterou jsou popsány vybrané metody dolování
43 dat a také je zde podrobně rozebrán algoritmus k -means. Tento algoritmus je využit pro
44 samotné dolování dat, které je zprostředkováno aplikací RapidMiner. Popis tohoto programu
45 a dalších vybraných nástrojů uzavírá třetí kapitolu.

46 Implementační část této práce je popsána kapitolou **čtvrtou**. Je členěna na oddíly,
47 které odpovídají vybraným částem architektury této práce. Elementy architektury, které
48 jsou popsány vlastní podkapitolou, jsou rozšířeny o zjednodušený návrh struktury v podobě
49 grafických schémat.

50 Kapitola **pátá**.

51 Kapitola **šestá**.

52 **Přílohy**.

53 Kapitola 2

54 XMPP

55 V následující kapitole jsou, pro usnadnění a jednodušší pochopení, rozebrány základní sta-
56 vební kameny protokolu Extensible Messaging and Presence Protocol (XMPP). Konkrétně
57 jsou zde popsány stávající vlastnosti implementace, architektura protokolu XMPP obecně
58 [22, 23] a další detaily protokolu [1, 24, 14]. Vzhledem k požadavkům na dolování v da-
59 tech popsaných v následující kapitole je kladen důraz na vybraná rozšíření [21, 12]. Tato
60 rozšíření tvoří základ pro některé rozšířené statusy, jako je například User Tune [18], User
61 Mood [20], User Location [6] a další. Další informace použité pro popis a pochopení XML
62 jazyka byly čerpány z [10, 9].

63 Vznik samotného protokolu XMPP je datován do roku 2004 (březen), kdy na něj byl
64 přejmenován Jabber. Původní projekt Jabber byl vytvořen roku 1998 autorem Jeremie
65 Millerem, který ho založil za účelem vytvořit svobodnou otevřenou IM službu. Uvedený
66 projekt měl obsahovat tři základní vlastnosti, do kterých se zahrnují jednoduchost a sro-
67 zumitelnost pro implementaci, jednoduchost v oblasti šíření a otevřenost podobě veřejně
68 dostupného popisu samotného protokolu. Základní vlastnosti a výhody klientů a serverů
69 budou podrobněji popsány níže. Roku 1999, 4.ledna byl vytvořen první server se jménem
70 Jabber. Komunita vývojářů se chopila iniciativy a vytvořila klienty, kteří dokázali se ser-
71 verem komunikovat, pro různé platformy (Linux, Macintosh, Windows). Roku 2004 byl
72 protokol XMPP přidán mezi RFC¹ dokumenty. Základní norma popisující obecnou struk-
73 turu protokolu je RFC 3920 [22] a RFC 3921 [23], který se zaměřuje na samotný instant
74 messaging a zobrazení stavu. Další zdokumentovaná rozšíření jsou vydávána v podobě tzv.
75 XEP (XMPP Extension Protocol) dokumentů, které jsou známé také pod starším názvem
76 JEP (Jabber Enhancement Proposal). Dnešní počet těchto norem se blíží k číslu 300. Každý
77 XEP obsahuje stav vývoje (schválení), ve kterém se zrovna nachází.

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

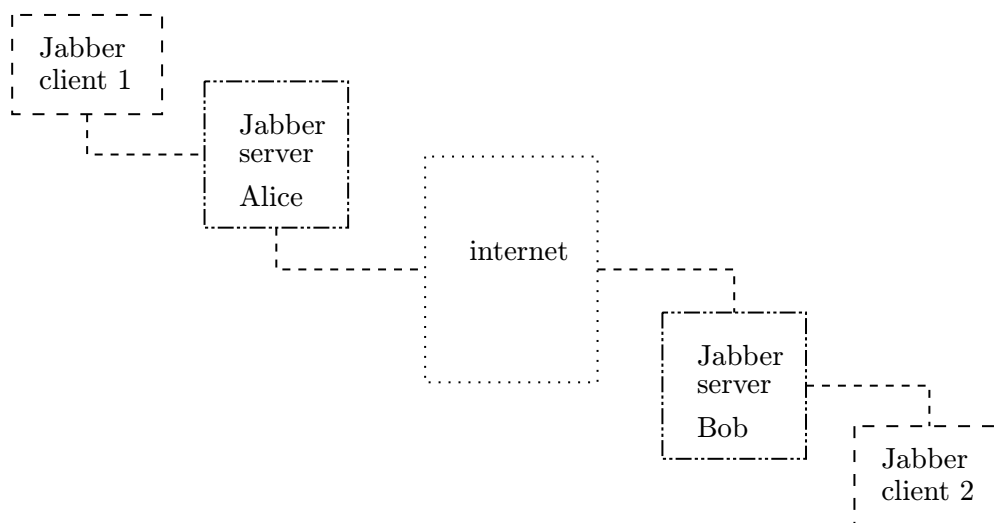
81 2.1 Architektura

82 Dobře navržená internetová technologie je tvořena správně fungujícími komponenty, které
83 mezi sebou dokáží vytvořit spojení a následně započít komunikaci. Pro popis Jabber ar-
84 chitektury v této práci bylo čerpáno z [1, 24]. Tato struktura se nejvíce podobá struktuře
85 posílání e-mailů. Hlavní předností Jabber sítě je, tak jako u elektronické pošty, její decentra-

¹RFC request of comments – žádost o komentáře

lize. 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.



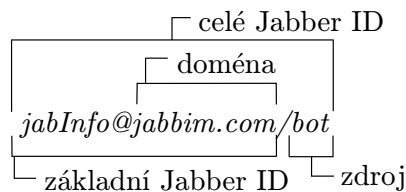
Obrázek 2.1: Distribuovaná architektura Jabber.

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“ JID², 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* [22]. 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.

²uživatelské jméno



Obrázek 2.2: Rozebraná struktura Jabber ID.

112 Další vymoženost JabberID oproti e-mailové adrese je jeho možnost používat prakticky
 113 libovolné národní znaky u doménových jmen a uživatelských účtů [24]. Využíváním kó-
 114 dování UNICODE, se XMPP stává plně mezinárodní a není jako jiné protokoly omezen
 115 rozsahem ASCII tabulky. Přestože je tato vymoženost k dispozici, doposud není žádným
 116 výrazným způsobem využívána.

117 Klient

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

- 122 1. komunikace s jedním Jabber serverem pomocí TCP socketu, který garantuje spolehlivé
 123 doručení zpráv na rozdíl od UDP. Nad tímto transportním protokolem dále běží
 124 kryptografický protokol TLS, který zabezpečuje komunikaci klient-server a server-
 125 server.
- 126 2. rozparsování a následná interpretace příchozí XML zprávy „stanza“ (kapitola 2.3)
- 127 3. porozumění sadě zpráv (*message*, *iq*, *presence*) z Jabber jádra [22]

128 Server

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

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

145 V tabulce 2.1 jsou shrnuty informace o serverech Jabberu. První sloupec tvoří jméno,
 146 následuje programovací jazyk, v němž je napsán. Většina aplikací pro servery je vydá-
 147 vána pod licencí GPL³. U všech aplikací byla zkoumána nejaktuálnější verze. Její číslo
 148 lze nalézt ve třetím sloupci. Všechny servery lze provozovat na operačním systému Li-
 149 nux a Windows. Na platformě Mac OS mohou být použity všechny zde jmenované vyjma
 150 jabberd2. Pět z šesti zde představených programů pro server Jabber jsou stále vyvíjeny,
 151 tedy kromě jabberd14. Hlavním účelem tabulky je prezentovat důležité vlastnosti serverů
 152 v oblasti podpory rozšířených statusů. Jedná se o standardy *pubsub*⁴ (XEP-0060) [12] a o
 153 jeho verzi zaměřenější více na uživatele *pep*⁵ (XEP-0163) [21]. Obě tato rozšíření tvoří ne-
 154 zbytnou základnu pro *rozšířené statusy* a proto je jejich podpora jak u serverů, tak klientů
 155 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
jabberd2	c	2.2.11	NE	NE
jabberd14	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.

156 Z výše uvedené tabulky je zřejmé, že aplikace pro servery, které jsou stále ve vývoji,
 157 podporují tzv. *rozšířené statusy*. Tedy kromě programu jabberd2.

158 2.2 XML

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

165 Ukázka možné struktury dokumentu psaného jazykem XML je zobrazena na příkladu
 166 2.1. Standardně je předpokládáno, že je psán v kódování UTF-8 [10], ale je-li jako v tomto
 167 případě použito jiné, musí být konkrétní kódování uvedeno na jeho počátku. V opačném
 168 případě nemusí být obsah správně zobrazen. Na začátku dokumentu se také uvádí verze
 169 XML, ve které je dokument psán (1. řádek příkladu). Následuje kořenový element, který
 170 je uzavřen na samotném konci dokumentu. 4. řádek prezentuje možnost použití prázdného
 171 elementu, který obsahuje jeden atribut s názvem zkratky fakulty. Velký význam zde mají
 172 úhlové závorky. Jsou jimi z obou stran obaleny všechny elementy.

³General Public License — všeobecná veřejná licence GNU

⁴Publish-Subscribe

⁵Personal 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 [24]. Zkráceně je ozna-

201 čována pomocí dvou počátečních písmen *Info/Query* neboli *IQ*. Na rozdíl od elementu
 202 *message* tvoří *iq* spolehlivější přenos, optimalizovaný pro výměnu dat (binární data). K
 203 dalším rozdílům patří povinnost příjemce odpovědět na každou přijatou zprávu, neboli po-
 204 tvrdit její doručení. Skutečnost, že je na právě požadovanou zprávu odpovězeno, zajišťuje
 205 parametr *id*. *Iq* dotaz nebo odpověď musí obsahovat stejnou hodnotu tohoto atributu jako
 206 zpráva vytvořená žádajícím subjektem. Další povinný atribut rozděluje *iq* na čtyři typy.
 207 Jednotlivé žádosti na proces nebo akci jsou posílány samostatně [23]. V příloze E je uve-
 208 dena rozsáhlejší struktura tohoto elementu. Použití nachází v případech, které nastavují,
 209 žádají nebo informace posílají. Tato struktura je využívána pro novou registraci, posílání
 210 seznamu kontaktů a další.

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

```

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

```

Příklad 2.3: Použití elementu *iq*.

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

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

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

```

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

```

Příklad 2.4: Použití elementu *presence*.

230
 231 Obsáhlejší struktura elementu *presence* je zobrazena v příloze E, kde je rovněž k nalezení
 232 přehled všech možných stavů.

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

244 2.4 Rozšíření

245 Dále se tato práce zabývá rozšířeními protokolu XMPP o další vlastnosti, k jejichž popisu
246 slouží XEP. Pro tuto práci jsou nepostradatelné „statusy“, pro které tvoří základ standardy
247 XEP-0060 [12] a XEP-0163 [21] zkráceně PEP⁶. Obě tato rozšíření umožňují strukturovaně
248 pracovat, používat a přenášet další XEP protokoly. Jako příklady relevantní k práci jsou zde
249 uvedeny protokoly *User Location* (kde se uživatel právě nachází) [6], *User Tune* (co uživatel
250 poslouchá za hudbu) [18], *User Mood* (aktuální nálada uživatele) [20] a *User Activity* (co
251 uživatel právě dělá) [11]. Jsou to tedy protokoly založené na PEP, které vyžadují podporu
252 nejen v klientech, ale i na straně serveru (zobrazuje tabulka 2.1). S touto informací úzce
253 souvisí další protokol XEP-0115 [7], který umožňuje zjistit podporované schopnosti klienta,
254 případně, které informace je ochoten přijímat. Tato vlastnost bude popsána níže v části
255 zabývající se podporovanými vlastnostmi.

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

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

```
1      <iq from='user@jabbim.com' type='set' id='publ'>
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.

⁶Personal Eventing via Pubsub

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

270 Podporované vlastnosti

271 Jednotlivá rozšíření protokolu XMPP jsou nepovinná, a proto nemusí být ve všech klient-
272 ských aplikacích podporována. Pro zjištění podporovaných rozšíření se používá XEP-0115
273 Entity Capabilities [7]. Toto rozšíření výrazně snižuje počet a velikost komunikací a přenosů
274 zpráv mezi uživateli. Dotazem zobrazeným na příkladu 2.6 je zjištěna schopnost jednotli-
275 vých klientů, kterou následně server využije pro správné směrování rozšířených statusů.
276 Všechny zde zmiňované rozšíření a protokoly z této kapitoly je možné u každého klienta
277 (seznam klientů obsahuje tabulka v příloze E) vyčíst z atributu *ver* (druhá část u atributu
278 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.

279 Další rozšíření

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

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

295 Díky základní podmínce XMPP protokolu (otevřenost) existuje mnoho různých aplikací,
296 pomocí kterých lze v síti Jabber komunikovat. S programy, používanými uživateli, úzce
297 souvisí další zde implementované rozšíření. Jedná se o realizaci *Software Version* dokumentu
298 [17], který se právě zabývá získáváním informací o samotných aplikacích. Je-li toto rozšíření
299 podporováno je díky němu možné zjistit jméno a verzi používané aplikace. Informace o
300 operačním systému často nejsou kvůli bezpečnosti ani vyplněny. Podrobnější informace o
301 softwarové výbavě klienta je možné zjistit pomocí XEP [7], o kterém již bylo dříve psáno v
302 odstavci zabývajícím se podporovanými vlastnostmi klientských aplikací.

303 S rozšířením tzv. „chytrých“ mobilních zařízení mezi širší veřejnost vzniklo několik no-
304 vých disciplín spojených s určováním zeměpisné polohy, jako je například geocaching. Geo-
305 grafická poloha je přenášena ve formě souřadnic popisující přímo zeměpisnou šířku a délku.
306 Současně lze informaci o poloze přenášet i slovně ve formě adresy. Příkladem slovního po-
307 pisu je ulice, číslo popisné, město a další. Mnoho aplikací, které mají k dispozici GPS
308 přijímač, vysílají a aktualizují zeměpisné informace automaticky, například po určité době
309 nebo změně polohy o určitou vzdálenost. Toto a další níže popsané rozšíření jsou postaveny
310 na již zmiňovaném PEP. Některé části protokolů jsou zjednodušeny a připraveny tím pro
311 „mobilní instant messaging“.

312 Pro sdělení informací o stavu klienta není v základní verzi Jabberu mnoho. Pomocí
313 presence je možné „pouze“ prozradit, zda je uživatel připraven komunikovat nebo je mo-
314 mentálně nedostupný a to v několika verzích lišících se délkou nepřítomnosti. Pokročilejší
315 nastavení statusu nabízí *User Mood* [20] a to ve formě sdělení současné nálady, jako je
316 například radost. Další možné upřesnění činnosti uživatele jsou definovány v *User Activity*
317 [11], kde každá činnost je složena z povinné obecné kategorie a nepovinné, která informaci
318 upřesňuje. Příkladem může být *eating* a *having_a_snack* tj. uživatel jí, uživatel svačí.

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

325 Podpora rozšíření v aplikacích je ukázána v tabulce v příloze G. Z této tabulky vyplývá,
326 že rozšířenost aplikací podporující výše popsaná rozšíření je poměrně malá. Například v
327 předcházející zmiňované části o poslouchané hudbě, při stavu, kdy program toto rozšíření
328 nepodporuje, je posíláno pomocí normální presence. Jméno skladatele, alba a další podrobn-
329 nosti jsou shrnuty do statusu, tudíž jsou doručeny všem uživatelům ze seznamu kontaktů.

330 Kapitola 3

331 Data mining

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

338 Na začátku kapitoly je rozebrán pojem získávání znalostí databází, jehož jednu pod-
339 statnou část tvoří samotný data mining. Dále je vysvětlena základní terminologie, pro kterou
340 bylo čerpáno z [8]. Cílem první podkapitoly je přiblížení způsobu, jakým byla data uložena
341 v databázi, připravena k samotnému data miningu. Celá druhá podkapitola je věnována vy-
342 braným metodám pro dolování dat a vlastnostem, které je od sebe navzájem odlišují. Jsou
343 zde rozebrány *asociační pravidla*, pro jejichž popis bylo čerpáno z [2]. Pro ostatní metody,
344 které jsou popsány dále, byla jako zdroj informací použita kniha [5]. Poté následuje třetí
345 podkapitola, která se podrobněji zabývá jednou z metod pro dolování dat a to *shlukováním*.
346 Obsahem této části jsou již konkrétní algoritmy pro shlukování dat [27, 3] a také metoda
347 *k-Means* využívaná v praktické části této práce. Kapitulu uzavírá stručný přehled vybra-
348 ných programů pro data mining a podrobnější seznámení s nástrojem *RapidMiner*, který je
349 v této práci využíván pro samotné dolování.

350 Terminologie

351 Pojem data mining neboli česky dolování dat se začal ve vědeckých kruzích objevovat počát-
352 kem 90. let 20. století. První zmínka pochází z konferencí věnovaných umělé inteligenci (IJ-
353 CAI'89¹—mezinárodní konference konaná v Detroitu, AAAI'91² a AAAI'93—americké kon-
354 ference v Californii a Washingtonu, D.C) [2].

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

- 361 1. počet záznamů neboli objektů v databázi

¹International Joint Conference on Artificial Intelligence

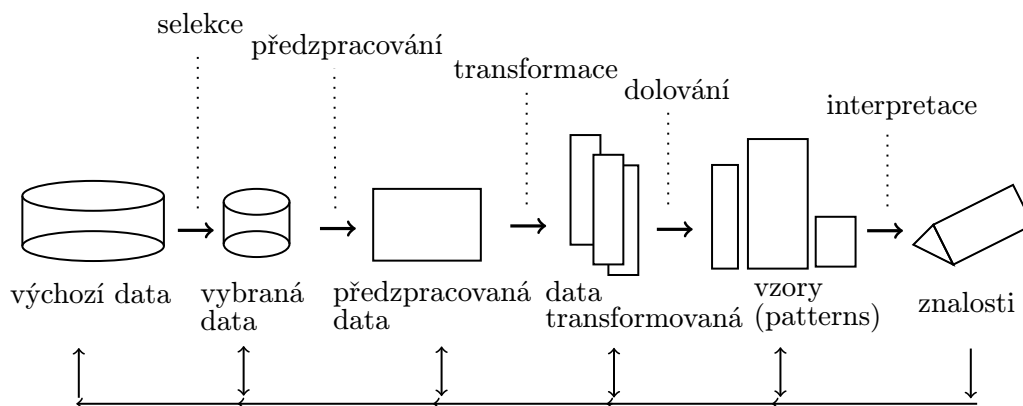
²Association for the Advancement of Artificial Intelligence

362 2. počet polí neboli atributů objektů v databázi

363 Proces data mining je pouze jedna část z odvětví nazývané dobývání znalostí z da-
364 tabází neboli KDD³ definované níže v definici 3.0.1. Vznik disciplíny KDD je důsledkem
365 nepřehledného množství automaticky sbíraných dat, která je potřeba dále využívat. Pod-
366 statným znakem celého procesu je správnost reprezentace výsledků formou, která má k
367 uživateli nejbližší. Jako příklad bude uvedena implikace ve tvaru rozhodovacích pravidel,
368 asociační pravidla, rozhodovací stromy, shluky podobných dat a další. Základem KDD je
369 praktická použitelnost metod. Očekává se zjištění nových skutečností namísto prezentování
370 již známých informací.

371 **Definice 3.0.1** KDD je chápáno jako interaktivní a iterativní proces tvořený kroky se-
372 lekce, předzpracováním, transformace, vlastního „dolování“ (data-mining viz 3.0.2) a in-
373 terpretace [2].

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



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

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

³Knowledge Discovery in Database

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

- 393 • **čištění dat** – fáze, ve které jsou nepodstatné údaje odstraněny z kolekce.
- 394 • **integrace dat** – kombinování heterogenních dat z několika zdrojů do společného
395 jediného zdroje.
- 396 • **výběr dat** – rozhodování o relevantních datech.
- 397 • **transformace dat** – také známý jako konsolidace dat. Fáze, ve které jsou vybraná
398 data transformována do formy vhodné pro dolování.
- 399 • **data mining** – zásadní krok, ve kterém jsou aplikovány vzory na data.
- 400 • **hodnocení modelů** – vzory dat zastupují získané znalosti.
- 401 • **prezentace znalostí** – konečná fáze, zjištěné poznatky jsou reprezentovány uživa-
402 teli. Tento základní krok využívá vizualizační techniky, které pomáhají uživa-
403 telům porozumět a správně interpretovat získané výsledky.

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

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

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

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

414 Mezi další čteně se vyskytující pojmy v tomto odvětví patří například data, znalosti a
415 informace. Tyto termíny jsou často mezi sebou zaměňovány, proto jsou níže jejich významy
416 striktně definovány tak jako v [8].

417 Jedna z několika existujících definic pojmu data je uvedena v definici 3.0.3, která je po-
418 pisuje z pohledu informačního. Data často nemají sémantiku (význam) a bývají zpracována
419 čistě formálně.

420 **Definice 3.0.3** Data jsou z hlediska počítačového pouze hodnoty různých datových typů.

421 Informace lze chápat jako data, která byla obohacena o sémantiku (význam), jsou tedy
422 již zpracovaná a interpretována uživatelem. Znalosti, jsou řazeny do stejné kategorie jako
423 informace, ale jejich interpretace bývá ještě složitější. Často bývají tvořeny shluky informací,
424 proto jsou reprezentovány jako odvozené informace. Podle studijní opory [8] jsou znalosti
425 informace, které jsou zařazeny do souvislostí.

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@jabbin.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 [22].

Temporální data

Druhá podstatná transformace, pro kterou bylo čerpáno z [25], 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.

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

459 K dalšímu popisu temporálních databází nyní budou charakterizovány tři důležité po-
 460 jmy, pomocí nichž jsou databáze dále děleny. Prvním pojmem je *granualita*, která udává
 461 nejmenší časovou jednotku, kterou databáze rozlišují. Hodnoty, které může nabývat, jsou
 462 hodina, den, rok a další. Velikost granuality ovlivňuje velikost objemu dat, který je přímo
 463 úměrný s přesností záznamů. Dalším pojmem je *čas platnosti*, která reprezentuje období,
 464 kdy je daný fakt v modelovém světě pravdivý. Posledním termínem je *čas transakce*, která
 465 definuje přítomnost faktu v databázi a možnost jej získat. Čas transakce a čas platnosti jsou
 466 na sobě nezávislé a definují dvě rozdílné časové osy, kdy každá může disponovat s jinou gra-
 467 nualitou. Souhrnný název pro výše uvedené tři pojmy, který se používá, je systém časových
 468 razítek. Při použití těchto systému lze následně tvořit dotazy zaměřené na různá časová
 469 období, jako je minulost, přítomnost a budoucnost. Tyto dotazy jsou velmi jednoduché a to
 470 díky rozšířenému jazyku TSQL, který vychází z klasické podoby dotazovacího jazyka SQL.

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

476 Konkrétní příklad možné transformace je ukázán na části tabulky *presence* 3.1, která
 477 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*.

478 butem *dateStart* vymezuje interval, kdy daná hodnota byla nebo je platná. Tato entita
 479 je pouze příkladem jak by mohla daná transformace vypadat. V této práci se žádná nová
 480 tabulka nevytvářela a ani se nemodifikovala již vytvořená. Celá transformace je popsána v
 481 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*.

482 Postup jednotlivých kroků při transformaci časových údajů z tabulky *presence*, je zob-
 483 razen v následujícím výčtu. Tento zjednodušený popis algoritmu vyžaduje dva vstupní
 484

485 parametry. První je JabberID uživatele, jehož položky v databázi mají být transformovány.
486 Druhá nutná položka je datum, které bude sloužit jako upřesňující vstupní interval.

Data z tabulky presence jsou transformována na vektor ϑ 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

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

487 kde $N = 288$.

- 488 1. Vstupním parametrem je datum, které vymezuje data pro transformaci. Toto datum je
489 převedeno na dvě data (počátek a konec dne), která tvoří hraniční body v intervalu ι .
- 490 2. Výběr dat z databáze, která splňují časové období definované intervalem ι a uživatel
491 ID . Uložení těchto dat do množiny Γ .
- 492 3. Pokud zadanému dotazu neodpovídá žádný řádek z tabulky, jsou data označena jako
493 prázdná. Výstupní transformovaný vektor ϑ je naplněn hodnotami reprezentujícími stav
494 *Unavailable*. Algoritmus je **ukončen**.
- 495 4. Zjištění prvního statusu toho dne.
 - 496 4.1 Převod data o den dřívejšího na interval ι_1 , například $\langle 2011-08-22\ 00:00:00, 2011-$
497 $08-22\ 23:59:59 \rangle$.
 - 498 4.2 Výběr dat vyhovujícím intervalu ι_1 a uživateli ID z databáze.
 - 499 4.3 Výběr posledního záznamu z množiny dat získaných z předešlého dotazu.
 - 500 4.4 Nalezení poslední presence a uložení její hodnoty do λ .
- 501 5. Výběr následujícího záznamu z množiny Γ .
- 502 6. Výpočet zda je časový interval mezi vybraným a následujícím záznamem větší jak δ .
 - 503 6.a Časový interval je větší než interval δ .
 - 504 6.1 Do výsledného vektoru ϑ je ukládána hodnota presence daného záznamu.
 - 505 6.2 Opakuj předešlý bod **6.1** kolikrát je interval δ menší než rozdíl mezi časy
506 vybraného a následujícího záznamu.
 - 507 6.b Časový interval je menší než interval δ .
 - 508 6.1 Jsou vybírány další záznamy z množiny Γ dokud rozdíl mezi časy v sekundách
509 není větší než interval δ .
 - 510 6.2 Jednotlivé presence záznamů jsou ukládány do pomocného pole, transformo-
511 vané do kvantitativních hodnot, jak je uvedeno v úvodu této podkapitoly.
 - 512 6.3 Každý prvek pole je vynásoben počtem sekund, zastoupených v daném inter-
513 valu.
 - 514 6.4 Výběr největšího prvku z pomocného pole a uložení jej do výsledného vektoru
515 ϑ .
- 516 7. Celý proces je opakován od bodu **5**, dokud množina Γ není prázdná.

517 3.2 Metody dolování dat

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

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

- 526 • schopnosti reprezentace shluků (např. otázka lineární separability)
- 527 • srozumitelnosti nalezených znalostí pro uživatele (symbolické vs. subsymbolické me-
528 tody)
- 529 • efektivnosti znovupoužití nalezených znalostí
- 530 • vhodnosti typů dat
- 531 • a další...

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

535 Při popisu asociačních pravidel, která jsou založena na syntaxi *IF-THEN*, bylo čerpáno
536 z [2]. Jejich rozšíření se datuje do 90. let 20. století, kdy byly panem Agrawalem před-
537 staveny v souvislosti s analýzou „nákupního košíku“. Použitelnost bude vysvětlena právě
538 na příkladu analýzy nákupního košíku. Podstata příkladu je tvořena zákazníkem a jeho
539 systémem nakupování. Jsou zjišťovány produkty, které jsou nakupovány současně. Hledají
540 se neboli jsou vytvářeny společné vazby (asociační pravidla) mezi výrobky a určuje se je-
541 jich spolehlivost. Na základě těchto závislostí je upravováno umístění jednotlivých výrobků.
542 Obecně jsou tedy asociační pravidla považována za konstrukci, která z hodnot jedné trans-
543 akce odvozuje možnost výskytu závislostí v jiných transakcích. Jsou tedy hledány všechny
544 vnitřní závislosti existující mezi daty.

545 K dalším metodám pro data minig patří klasifikace, která bude opět vysvětlena na pří-
546 kladu, převzatého z [8]. Podle obsahu databáze nebo dotazníku bude každý klient banky
547 zařazen do různých krizových skupin. Na základě těchto skupin pracuje „credit skóring“,
548 jež klientovi poskytne nebo odepře například úvěr v bance. Další příklady využití jsou na-
549 příklad ve zdravotnictví. Na základě zdravotního stavu pacienta a jeho příznaků, je pacient
550 zařazen do tříd, které reprezentují jednotlivé nemoci. Klasifikací jsou, podle [5], jednotlivé
551 zkoumané elementy rozděleny (podle hodnot atributů) do vhodných kategorií, které jsou
552 předem vytvořeny z navzájem podobných objektů (tvorba profilů třídy). Při této metodě je
553 upřednostňována přesnost před jednoduchostí a rychlostí. Zdroje klasifikovaných objektů
554 jsou většinou tvořeny jednotlivými řádky v databázi. Vzory dat vytváří instance, jejichž
555 vlastnosti reprezentují atributy vyjádřené číselnou hodnotou.

556 Na modelech je založen třetí typ metod pro dolování dat. Základem modelů jsou tré-
557 novací data. Dále uvedené příklady vybraných klasifikačních modelů, byly čerpány z [5].
558 Především jsou to rozhodovací stromy, neuronové sítě, statistické metody, klasifikační pra-
559 vidla a další.

Metoda, která je postavena na myšlence, kde chronologicky seřazená data a vývoj jejich hodnot v minulosti tvoří základ pro určení hodnot budoucích, se nazývá predikce. Je řazena mezi velmi známé procesy, které na základě získaných znalostí předpovídají následující vývoj. Předpokládá se, že na základě informací získaných z dat v minulosti, bude možné postavit modely, které se budou chovat stejně nebo alespoň podobně i v budoucnu. Využití naleznu v předpovědi počasí (z naměřených meteorologických hodnot se určují budoucí předpokládané teploty), při vývoji cen na burze a dalších. Podklady pro popis predikce byly čerpány z [2, 5].

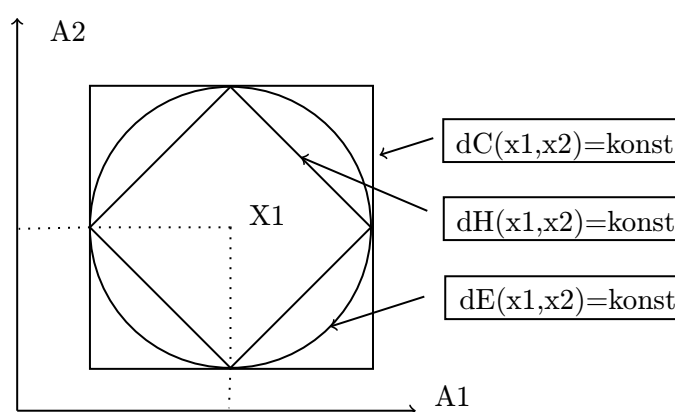
Poslední zde uvedená metoda je zaměřená na dělení objektů do předem neznámých skupin. Proces dělení probíhá na základě specifikace objektů a jejich odlišnosti od ostatních shluků. Tato část, pro kterou bylo čerpáno z [27, 3], bude podrobně rozebrána v následující podkapitole.

3.3 Shlukování

V této podkapitole je shlukování 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° 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

590 data k určitému matematickému modelu na základě společných optimalizovaných vlastností.
591 Většina procesů je založena na předpokladu, že jsou data generována pomocí standardních
592 statistik. Mezi zástupné metody této shlukovací analýzy se řadí Expectation–Maximization
593 (EM) a Self Organizing Oscillator Network, dále jen SOON. Algoritmus SOON je založen na
594 neuronové síti. Je to metoda vycházející z algoritmu SOM⁴ [27]. Metoda EM je rozšířením
595 algoritmu *k-means*, který bude podrobně rozebrán v následující části.

596 Hlavní princip *metody hierarchického shlukování* je založen na tvorbě stromové hierar-
597 chie shluků, která je známá pod názvem *dendrogram*. Hierarchické metody, podle [5], mohou
598 být rozděleny do dvou skupin a to na základě principu, kterým jsou dendrogramy vytvářeny.
599 První možnost je *aglomerativní přístup*, který shlukuje menší shluky, kdy výsledkem je jen
600 jeden. Druhý přístup, *divizní*, je založen na opačném předpokladu. Na počátku je tedy jeden
601 velký shluk, který je postupně rozdělován, dokud není počet shluků roven počtu objektů
602 [27]. Mezi zástupce této metody například patří algoritmus AGNES⁵.

603 K dalším metodám patří *metody založené na mřížce*, které kvantují datový prostor do
604 konečného počtu pravoúhlých buněk. Tyto buňky jsou uspořádány do víceúrovňové mříž-
605 kové struktury. Zmíněná struktura tvoří základ pro shlukové operace. Hlavní výhoda tohoto
606 přístupu je rychlost zpracování, které většinou nebere ohled na počet datových objektů. Čas
607 zpracování závisí pouze na počtu buněk v každé dimenzi kvantovaného prostoru. Mezi zá-
608 stupce metod založených na mřížce patří metoda STING — STatistical INformation Grid,
609 který pracuje se statickými informacemi uloženými v buňkách mřížky. Algoritmus je roz-
610 dělen do dvou částí. První si klade za cíl rekurzivně rozdělit datový prostor na pravoúhlé
611 buňky. Druhá fáze testuje spojitost mezi sousedy relevantních buněk [27]. Mezi další metody
612 založené na mřížce patří WaveCluster⁶, využívající vlnkové transformace k rozdělení pro-
613 storu dat. Tato transformace zdůrazňuje shluky v prostoru a objekty jim vzdálené potlačuje
614 [5].

615 *Metody založené na hustotě* vychází z *m-rozměrného* prostoru, ve kterém jsou zobrazeny
616 objekty ve formě bodů. Místa v prostoru s větší koncentrací objektů ve srovnání s ostatními
617 oblastmi jsou nazývány shluky. Výchozí předpoklad je existence okolí jednotlivých bodů
618 (sousedství). Jedna z charakteristik metod založených na hustotě je schopnost vypořádat
619 se s vzdálenými hodnotami, označovanými jako šum [27]. Jako příklad je uvedena metoda
620 DBSCAN⁷, která je založena na hustotě objektů v prostoru. U jednotlivých objektů je
621 zkoumáno jejich okolí. Algoritmus je ovlivňován dvěma parametry, velikostí shluku ϵ a
622 minimálním počtem objektů v daném shluku *MinPts*, které spolu úzce souvisí (viz [5]).
623 Bod splňující obě podmínky je označen za jádro. Za pomoci jader je rozšiřována množina
624 objektů spojených na základě hustoty. Obsahuje-li jádro x_1 ve svém okolí další jádro x_2
625 znamená to, že jádro x_1 je přímo dosažitelné z jádra x_2 . Tímto způsobem jsou vytvářeny
626 výsledné *shluky*. V opačném případě, body, které nesplňují dvě zmíněné podmínky, jsou
627 označeny jako *šum*.

628 Všechny zde doposud zmiňované metody poskytují dobré výsledky pouze s malým poč-
629 tem dimenzí, tak jak je to popsáno v [8]. S narůstajícím počtem atributů roste počet
630 nerelevantních dimenzí určených pro shlukování. S tímto také přibývá zvětšená produkce
631 zašumění a znesnadnění nalezení relevantních shluků. Data jsou roztroušena do mnoha di-
632 menzí a tím odpadá možnost použití vzdálenostních funkcí. Zmíněné problémy shlukování
633 velkých dat řeší dvě techniky *metoda transformace rysů* a *metoda výběru atributů*. Pro

⁴Self–Organizing Map

⁵AGglomerative Nesting

⁶Clustering Using Wavelet Transformation

⁷Density–Based Spatial Clustering of Applications with Noise

634 efektivní shlukování je možné použít například algoritmus CLIQUE⁸.

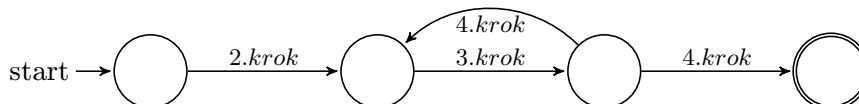
635 Metody rozkladu

636 Metody rozkladu rozdělují datové prvky do několika podmnožin, nazývané shluky. Počet
637 shluků musí být znám před zahájením samotného procesu. Přiřazení do konkrétních tříd
638 je, podle [27], jednoznačné nebo probíhá na základě míry příslušnosti objektů do shluků.
639 Pro velký počet objektů, se kterými se pracuje, jsou využívány různé iterační optimalizace.

640 Hlavním zástupcem u uvedených metod je algoritmus k -means, který je popsán níže.
641 Tvoří základ pro většinu metod shlukování nejen pro metody rozkladu. K dalším metodám
642 se řadí k -medoidů, k -modů, k -histogramů, fuzzy shluková analýza a další.

643 k -means

644 Shlukování pomocí algoritmu k -means je používáno pro data obsahující kvantitativní pro-
645 měnné a pro data, která nejsou příliš zašumělá. Základní proces je tvořen iterativním roz-
646 dělováním objektů do tříd na základě vzdáleností od jejich středů. Střed neboli centroid
647 shluku je vektor, jehož vzdálenost od součtu vzdáleností objektů v této třídě je minimální.
648 Celý tento proces je prezentován na obrázku 3.3 pomocí jednoduchého schématu konečného
649 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.

650 a středem je použita euklidovská vzdálenost⁹, která je vyobrazena na obrázku 3.2.

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

661 K -menas je algoritmus, kterým jsou přiřazovány objekty (vektory) x_n , kde $n = 1, \dots, N$,
662 do S_k , kde $k = 1, \dots, K$, shluků. V prvním kroku jsou určeny počáteční středy tříd, do
663 kterých se budou objekty shlukovat. Určení počátečních centroidů c_k probíhá například
664 náhodným výběrem K objektů nebo K prvních objektů souboru. Druhým krokem jsou
665 zkoumány jednotlivé vzdálenosti objektů x_n od počátečních středů c_j pomocí euklidov-
666 ské vzdálenosti. Na základě nejmenší zjištěné vzdálenosti mezi objektem a centroidem je
667 objekt zařazen do shluku, kterému náleží právě tento střed. Ve třetím kroku, tak jako u
668 kroku prvního, jsou hledány nové středy shluků. Nyní již však nejsou zvoleny náhodně, ale

⁸CLustering In QUEst

⁹mean = střed, centroid je vektor průměrů

669 spočítány. Jsou vypočítány na základě průměrných jednotlivých hodnot objektů a uložen
670 jako m -rozměrný vektor. Čtvrtým krokem se algoritmus dostává do konečné fáze, kdy mo-
671 hou nastat dva možné případy. Nově nalezené středy nejsou příliš vzdáleny od předchozích
672 centroidů a proto je algoritmus ukončen. Druhá častěji se vyskytující možnost iterativně
673 provádí algoritmus od druhého kroku, dokud neplatí první možnost nebo dokud se objekty
674 nepřestanou přemísťovat úplně. Při popisu tohoto algoritmu bylo čerpáno z [27, 2]. Níže
675 zobrazený algoritmus 3.1 prezentuje krok po kroku metodu k -means.

-
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].

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

680 3.4 Programy

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

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

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

701 Na vybraných technicky zaměřených vysokých školách existují skupiny, které se zabývají
702 výzkumem a vývojem nástrojů pro data mining. Jako příklad lze uvést již dříve zmiňovaný
703 FIT-miner z fakulty informačních technologií v Brně nebo také projekt LISp-Miner z Vy-

704 soké školy ekonomické v Praze. LISp-Miner je otevřený akademický systém určený pro
705 výuku a výzkum metod pro dobývání znalostí z databází.

706 **RapidMiner**

707 RapidMiner je, tak jak je popsán na oficiálních stránkách produktu [26], celosvětově nej-
708 používanější open-source systém pro dolování dat. Je možné jej používat jako samotnou
709 aplikaci nebo jej začlenit jako komponentu do vlastních výrobků v podobě knihovny pro ja-
710 zyk Java. Pro zájemce je nabízen také ve verzích pro firmy, které jsou rozdílné v poplatcích,
711 podpoře pro zákazníka, záruce a dalších balíčků služeb zajišťující celkovou komplexnost a
712 spolehlivost produktu.

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

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

730 Kapitola 4

731 Implementace

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

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

742 Třetí část této kapitoly je zaměřena na Jabber klienta neboli robota, který je imple-
743 mentován v jazyce C++ pouze s konzolovým rozhraním. Robot v této práci hraje roli
744 pasivního uživatele, který informace pouze přijímá. Ve vybraných případech dokáže uží-
745 vatele ze svého seznamu kontaktů vyzvat k zaslání odpovědi s informacemi na odpověď,
746 o kterou žádal. Struktura samotného robota je popsána níže a reprezentována obrázkem
747 4.3. V části o Jabber robotovi jsou také uvedeny informace o knihovně *gloox*, kterou je
748 zprostředkovány všechny náležitosti Jabber komunikace.

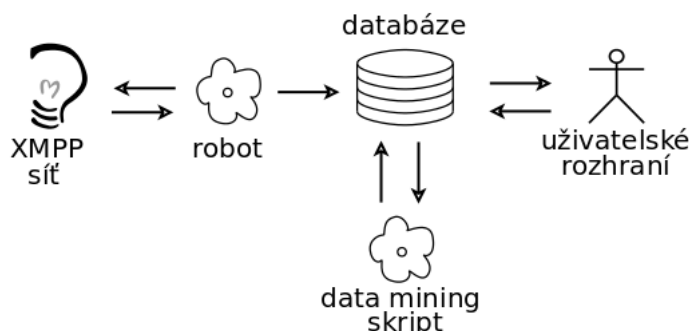
749 4.1 Architektura

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

755 První část schématu, která je nazvaná XMPP síť a je na obrázku 4.1, je obecně popsána
756 v druhé kapitole. Vztah a vzájemná komunikace s robotem probíhá pomocí internetové sítě,
757 kdy XMPP síť souhrnně reprezentuje jednotlivé prvky, jako jsou Jabber servery, uživatelé
758 a jejich klienty. Jak je patrné, výměna informací mezi těmito částmi probíhá obousměrně.
759 Druhý element schématu, *robot*, je charakterizován v podkapitole níže, kde je podrobně
760 popsán návrh jeho struktury a vyzdviženy jeho základní vlastnosti a schopnosti. Dalším
761 blokem je *databáze*, která reprezentuje datové úložiště. Do něhož jsou za pomoci jedno-
762 směrného kanálu ukládána data, která jsou získávána z toku informací proudících mezi
763 robotem a XMPP sítí. Přehled vybraných jednotlivých tabulek a jejich atributů je možné

764 nalézt v následující kapitole, která se zabývá návrhem databáze.

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



Obrázek 4.1: Struktura architektury bakalářské práce.

770 Poslední částí je *uživatelské rozhraní*, které je implementováno pomocí webových služeb.
771 Jako příklad lze uvést službu, která zobrazuje status uživatele na webové stránce. Uživatel
772 pouze musí vlastnit Jabber účet a mít přidáného tohoto robota do seznamu kontaktů.

773 4.2 Databáze

774 Data, která jsou sbírána robotem, jsou ukládána do objektově-relační databáze Postgre-
775 SQL¹. PostgreSQL neboli také Postgres byl použit ve verzi 8.4.7 a je provozován na opera-
776 čním systému Ubuntu.

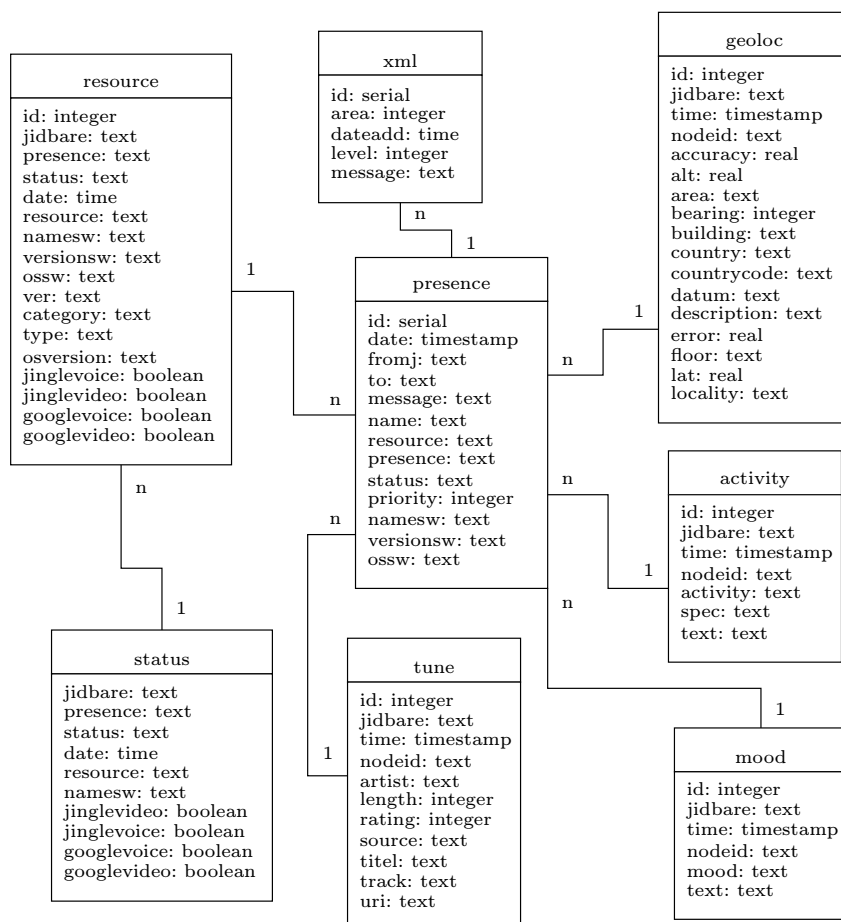
777 Návrh databáze

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

782 V době návrhu databáze nebylo zcela zřejmé, která data budou následně analyzována.
783 Z tohoto důvodu se struktura databáze snaží zachytit všechna „důležitá“ data. Za tímto
784 účelem je v návrhu databáze obsažena tabulka *xml*, která je nositelem obsahu jak všech
785 přijatých, tak i odeslaných zpráv. Tabulky *debug*, *level* a *logarea*, které v zúženém návrhu
786 databáze, uvedeném na obrázku 4.2, nejsou zobrazeny, jsou určeny pouze jako doplňkové
787 informace k typu zprávy. Tabulku *xml* je možné nahradit jednotlivými dalšími tabulkami,
788 které jsou zaměřeny na konkrétní data. Příkladem entity, která již obsahuje konkrétní data
789 bez *xml* prvků, je tabulka *message*, která je taktéž vyobrazena v příloze E. Jejím obsahem
790 jsou zprávy vzniklé při Jabber komunikaci mezi uživatelem a robotem, jehož schopnosti
791 budou popsány níže.

¹vyvinul se z projektu Ingres

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



Obrázek 4.2: Vybraná část struktury databáze.

801 Většina rozšíření standardu XMPP, která jsou popsána ve druhé kapitole, jsou pre-
802 zentována pomocí pěti tabulek. Konkrétně to jsou tabulky *geoloc*, *tune*, *mood*, *activity* a
803 *vcard*. Pro obsah a názvy atributů všech pěti tabulek s rozšířeními se staly vzory dokumenty
804 jednotlivých XEP, které je definují. Tabulky kromě těchto atributů obsahují také položku
805 *jidbare*, která, jak již název naznačuje, obsahuje pouze čisté JabberID bez resources. Tato
806 skutečnost vyplývá již ze samotného návrhu XEP protokolů. V tabulce *vcard* jsou také
807 položky, jejichž obsah je tvořen fotografiemi. Obrázky jsou zde uloženy ve stavu, tak jak
808 jsou přenášeny po síti, tedy pomocí datového formátu Base64. Base64 je algoritmus, který
809 převádí binární data na řetězec, který je tvořen pouze znaky ASCII tabulky.

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

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

824 4.3 Robot

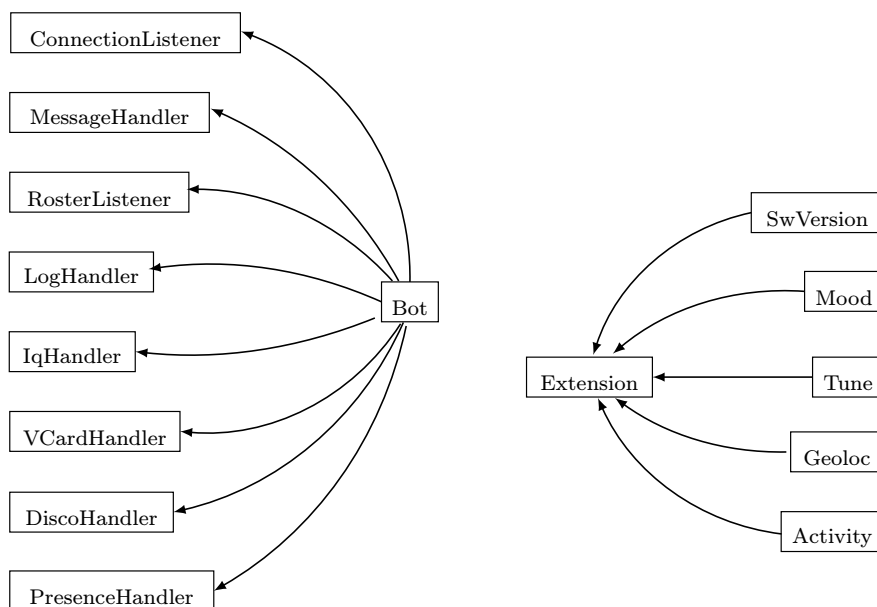
825 V této části bude popsána aplikace, která zajišťuje real-time komunikaci s ostatními en-
826 titami Jabber sítě. Konzolový robot, který je vyvíjen pro operační systém Linux, také
827 vhodným způsobem pracuje s databází, kam jsou ukládána jednotlivá získaná data. Jabber
828 jádro síťové komunikace [22], je implementováno pomocí volně dostupné knihovny *gloox*. V
829 porovnání s jinými knihovnami disponuje lepší podporou a dokumentací. Gloox plně im-
830 plementuje standart XMPP Core [22] a z větší části i standard XMPP IM [23]. Dodatečně
831 je podporováno kolem třiceti XEP standardů mezi něž například patří vcard-temp [16] a
832 další.

833 Při implementaci robota byly využity základní prvky objektově-orientovaného progra-
834 mování. Každá třída zde zastupuje jednoznačně oddělitelné elementy robota. Ať už se jedná
835 o blok komunikující s databází, různá rozšíření nebo o jádro robota. Všechny zde jmenované
836 objekty hrají důležitou roli jak v části dekompozice na menší ucelené celky tak i v celkové
837 komplexnosti programu. Jako příklad lze uvést třídy, které reprezentují jednotlivá pro tuto
838 práci vhodná rozšíření, jejichž základ je v podobě XEP dokumentů.

839 Návrh robota

840 Struktura návrhu tříd robota, která je v základní části prezentována obrázkem 4.3, je velmi
841 podobná struktuře návrhu databáze, obrázek 4.2. Mnoho atributů z robota a z databáze jsou
842 ve vztahu jedna ku jedné. Samotná komunikace a řízení mezi těmito entitami je zprostředko-
843 váno pomocí knihovny *libpq*. Libpq je programátorské rozhraní pro databázi PostgreSQL. Je
844 to sada funkcí v jazyce C, které umožňují klientským programům zadávat dotazy na server
845 a následně na ně obdržet odpovědi. Pomocí této knihovny je řízen přenos dat mezi robotem,
846 který sbírá data, a databází, kam jsou ukládána. Navázání spojení a následné jednotlivé
847 dotazy probíhají v blokujícím režimu. V souboru *connect.h* jsou, v podobě předdefinova-
848 ných konstant, uvedeny všechny údaje nutné k navázání spojení ať už s databází nebo s
849 komunikačním Jabber serverem. Tudíž se stává hlavní částí určenou k editaci připojení a
850 Jabber účtu na němž robot běží. K dalšímu pouze textovému dokumentu patří *const.h*,
851 který je zaměřen k definování jednotlivých příkazů a dotazů pro databázi. Jsou to příkazy
852 pro tvorbu a aktualizaci tabulek, které jsou taktéž definované jako konstanty. Identifikační
853 názvy, pro snadné odlišení, jsou psány velkými písmeny a začínají textem *DB_*.

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



Obrázek 4.3: Vybraná část struktury robota.

859 Pro rozšíření podle dokumentů XEP popsaných v kapitole druhé, v části zabývající se
 860 implementovanými rozšířeními, je základ nadtřída Extension. Tato třída obsahuje základní
 861 funkce a parametry pro všechna rozšíření. Jako příklad lze uvést uživatelské jméno klienta,
 862 který dané rozšíření „šíří“ po síti. Třídy, které z Extension dědí a jsou prezentovány na
 863 návrhu robota uvedeném na obrázku 4.3, jsou následující: *Geoloc*, *Tune*, *Mood*, *Activity* a
 864 *SwVersion*. Samotnou strukturou těchto rozšíření, jak již bylo uvedeno výše, jsou pomocí
 865 parametrů kopírovány vlastnosti jednotlivých XEP dokumentů. Tedy ke každému atributu
 866 xml zprávy existuje jak proměnná uchováající její obsah tak i tzv. „get“ a „set“ metoda.
 867 Pomocí další metody jsou jednotlivé části rozšíření rozparsovány, uloženy do tříd a následně
 868 vloženy do databázi. Poslední doposud nezmiňované rozšíření v podobě *vcard*, je využito z
 869 knihovny gloox, která jej implementuje.

870 Základní část, kterou je možné považovat za jádro aplikace, je třída *Bot*. Tato třída je
 871 založena na již zmiňované C++ knihovně gloox. Je implementována za pomoci vybraných
 872 tříd, ze kterých dědí. Jsou to především „handler“ třídy, kterými je zajišťován přístup k
 873 jednotlivým zprávám. Nebo-li jejich prostřednictvím je získán přístup k samotným elemen-
 874 tům stanzy, jako je message, iq, a presence, jejichž vlastnosti a využití jsou popsány v druhé
 875 kapitole. Konkrétní seznam tříd, ze kterých je děděno, zobrazuje návrh robota na obrázku
 876 4.3. Jako příklad lze uvést třídu, z prostoru jmen gloox, *VcardHandler*, díky níž je možné
 877 získat a zpracovávat vcard uživatelů, kteří jsou v seznamu kontaktů robota. Samotné přidá-
 878 vání uživatelů do seznamu kontaktů je prováděno automaticky. Pouze na straně klienta je
 879 nutné požádat robota o autorizaci, která však proběhne automaticky. Jednotlivé kontakty
 880 nejsou žádným způsobem tříděny do skupin, tudíž existuje pouze jedna.

881 Kapitola 5

882 Pokračování práce

883 mozne pokracovani prace bych dal jako samostatnou kapitolu nebo potkapitolu (je nutne o
884 tom napsat aspon 1.5 stranky textu, protoze je to jeden bod zadani)

885 POKACOVANI =tridit lidi do podskupin v kotakt listu, naprikad podle zpravy za-
886 slene s dotazem k autorizaci=: kazda skupina by mohla mit jin prava, jine dotazy an
887 sluzby...mozne rozsirit o ruzne sluzby, ktere poskytuji server, podle dotazu....

888 Kapitola 6

889 Vyhodnocení výsledků

890 6.1 Manuální rozbor dat

891 –charakter dat, jaka data jsem nasbiral –popst jak jsem data prochael rucne, tedy manualni
892 dataming, pocet zazanmu, delka behu programu, velikost datatabze, pocet uzivatelu...

⁸⁹³ Kapitola 7

⁸⁹⁴ Závěr

⁸⁹⁵

Literatura

- [1] Adams, D.: *Programming jabber*. Sebastopol: O'Reilly, první vydání, 2002, 455 s.,
ISBN 05-960-0202-5.
- [2] Berka, P.: *Dobývání znalostí z databází*. Praha: Academia, první vydání, 2003, 366 s.,
ISBN 80-200-1062-9.
- [3] Bramer, M.: *Principles of Data mining*. London: Springer, první vydání, 2007, 343 s.,
ISBN 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., ISBN 02-625-6097-6.
- [5] Han, J.; Kamber, M.: *Data mining : concepts and techniques*. San Francisco: Morgan
Kaufmann Publisher, druhé vydání, 2006, 770 s., ISBN 15-586-0901-6.
- [6] Hildebrand, J.; Saint-Andre, P.: XEP-0080: User Location. [online], 15-09-2009, [cit.
6. 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. 6. 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.
6. 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., ISBN 80-716-9860-1.
- [11] Meijer, R.; Saint-Andre, P.: XEP-0108: User Activity. [online], 29-10-2008, [cit.
6. 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. 6. května 2011].
URL <http://xmpp.org/extensions/xep-0060.html>

- 926 [13] Mizzi, S.; Saint-Andre, P.: XEP-0292: vCard4 Over XMPP. [online], 02-26-2008, [cit.
927 6. května 2011].
928 URL <http://xmpp.org/extensions/xep-0292.html>
- 929 [14] Moore, D.; Wright, W.: *Jabber developer's handbook*. Indianapolis: Sams Publishing,
930 první vydání, 2004, 487 s., iISBN 06-723-2536-5.
- 931 [15] Nemrava, M.; Pospíšil, J.: Dolování dat a jeho aplikace. [online], 2006, [cit. 6. května
932 2011].
933 URL http://www.spatial.cs.umn.edu/paper_ps/dmchap.pdf
- 934 [16] Saint-Andre, P.: XEP-0054: vcard-temp. [online], 07-16-2008, [cit. 6. května 2011].
935 URL <http://xmpp.org/extensions/xep-0054.html>
- 936 [17] Saint-Andre, P.: XEP-0092: Software Version. [online], 02-15-2007, [cit. 6. května
937 2011].
938 URL <http://xmpp.org/extensions/xep-0092.html>
- 939 [18] Saint-Andre, P.: XEP-0118: User Tune. [online], 30-01-2008, [cit. 6. května 2011].
940 URL <http://xmpp.org/extensions/xep-0118.html>
- 941 [19] Saint-Andre, P.: XEP-0153: vCard-Based Avatars. [online], 16-08-2006, [cit. 6. května
942 2011].
943 URL <http://xmpp.org/extensions/xep-0153.html>
- 944 [20] Saint-Andre, P.; Meijer, R.: XEP-0107: User Mood. [online], 29-10-2008, [cit.
945 6. května 2011].
946 URL <http://xmpp.org/extensions/xep-0107.html>
- 947 [21] Saint-Andre, P.; Smith, K.: XEP-0163: Personal Eventing Protocol. [online],
948 12-07-2010, [cit. 6. května 2011].
949 URL <http://xmpp.org/extensions/xep-0163.html>
- 950 [22] Saint-André, P.: Extensible Messaging and Presence Protocol (XMPP): Core.
951 [online], 10-2004, [cit. 6. května 2011].
952 URL <http://tools.ietf.org/html/rfc3920>
- 953 [23] Saint-André, P.: Extensible Messaging and Presence Protocol (XMPP): Instant
954 Messaging and Presence. [online], 10-2004, [cit. 6. května 2011].
955 URL <http://tools.ietf.org/html/rfc3921>
- 956 [24] Saint-André, P.; Smith, K.; Troncon, R.: *XMPP : the definitive guide : building
957 real-time applications with jabber technologies*. Sebastopol: O'Reilly, první vydání,
958 2009, 287 s., iISBN 978-059-6521-264.
- 959 [25] WWW Stránky: Database Systems. [online], 2007, [cit. 6. května 2011].
960 URL
961 http://www.cs.uct.ac.za/mit_notes_devel/Database/Lates%t/index.html
- 962 [26] WWW Stránky: RapidMiner. [online], 2011, [cit. 6. května 2011].
963 URL <http://rapid-i.com/content/view/181/190/>
- 964 [27] Řezánková, H.; Húsek, D.; Snášel, V.: *Shluková analýza dat*. Praha: Professional
965 Publishing, druhé vydání, 2009, 218 s., iISBN 978-808-6946-818.

⁹⁶⁶ **Příloha A**

⁹⁶⁷ **Obsah CD**

⁹⁶⁸ **Příloha B**

⁹⁶⁹ **Manual**

⁹⁷⁰ **Příloha C**

⁹⁷¹ **Konfigurační soubor**

972 Příloha D

973 Slovník zkratek

974 **DNS** — Domain Name System

975 **GPG** — dkshckdsjvlsdjvodsvjdfokj

976 **IM služby** — dkshckdsjvlsdjvodsvjdfokj

977 **IP** — Internet Protocol

978 **JEP** — dkshckdsjvlsdjvodsvjdfokj

979 **JID** — dkshckdsjvlsdjvodsvjdfokj

980 **SASL** — dkshckdsjvlsdjvodsvjdfokj

981 **TCP** — dkshckdsjvlsdjvodsvjdfokj

982 **TLS** — dkshckdsjvlsdjvodsvjdfokj

983 **WWW** — dkshckdsjvlsdjvodsvjdfokj

984 **XEP** — dkshckdsjvlsdjvodsvjdfokj

985 **XML** — dkshckdsjvlsdjvodsvjdfokj

986 **XMPP** — dkshckdsjvlsdjvodsvjdfokj

987 **e-mail** — dkshckdsjvlsdjvodsvjdfokj

988 **jabber** — dkshckdsjvlsdjvodsvjdfokj

989 **klient** — dkshckdsjvlsdjvodsvjdfokj

990 **presence** — dkshckdsjvlsdjvodsvjdfokj

991 **server** — dkshckdsjvlsdjvodsvjdfokj

992 **stanza** — dkshckdsjvlsdjvodsvjdfokj

993 **vCard** — dkshckdsjvlsdjvodsvjdfokj

994 **ID3v1** — dkshckdsjvlsdjvodsvjdfokj

995 **ID3v2** — dkshckdsjvlsdjvodsvjdfokj

996 **ASCII** — dkshckdsjvlsdjvodsvjdfokj

997 **utf-8** — dkshckdsjvlsdjvodsvjdfokj

998 **Base64** — dkshckdsjvlsdjvodsvjdfokj

999 Příloha E

1000 Stanza - základní schéma

1001 Přehled základních elementů, které jsou využívány při Jabber komunikaci. Struktura jed-
1002 notlivých částí stanzy ukazuje pouze prvky relativní k této práci. Pomocí hranatých závorek
1003 je znázorněna množina, ze které musí být vybrán právě jeden prvek. Na místě uvozovek se
1004 očekává jakákoliv povolená hodnota.

1005 Iq

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

Algoritmus E.1: Popis elementu *iq*.

1006 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>
```

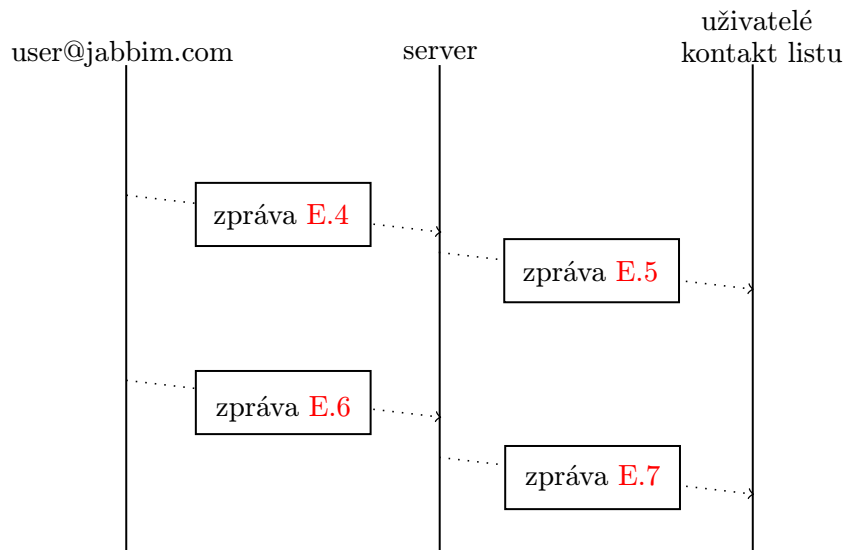
Algoritmus E.2: Popis elementu *message*.

```
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 E.3: Popis elementu *presence*.

1008 Přehled průběhu rozšíření

1009 Ukázka celého příkladu šíření statusu pomocí rozšíření *User Tune*. Uživatel *user* poslouchá hudbu a informuje server zasláním zprávy zobrazené v příkladu E.4.



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

1010

```

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            <source>Nigredo</source>
9            <title>1968</title>
10           <track>5</track>
11          </tune>
12        </item>
13      </publish>
14    </pubsub>
15  </iq>

```

Algoritmus E.4: Informování serveru o právě přehrávající hudbě.

1011 Server obdrží informace od klienta *user* zprávu o přehrávací hudbě. Pomocí elementu
1012 *message* ji přepoše všem uživatelům z kontakt listu uživatele *user*, kteří jsou pro odběr
těchto typů zpráv zaregistrováni. Tato struktura zprávy je prezentována na příkladu E.5.

```
1  <message from="user@jabbim.com" type="set"
2      to="jabinfo@jabbim.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 E.5: Server informuje uživatele podporující rozšíření o stavu *user@jabbim.com*.

1013 Zpráva o přehrávané hudbě je také přeposlána všem otevřeným spojením uživatele *user*,
1014 ukázáno na příkladě E.6.
1015

```
1  <message from="user@jabbim.com" type="set"
2      to="user@jabbim.com/doma" 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 E.6: Server přepoše informace o přehrávané hudbě všem otevřeným spojením
uživatele *user@jabbim.com*.

1016 Přestane-li uživatel *user* poslouchat/vysílat informace o přehrávané hudbě, provede to
1017 pomocí zprávy ukázané na příkladu E.7. Zpráva typu *iq*, ve které je položka *tune* nesoucí
informace o skladbě prázdná.

```
1 <iq from="user@jabbim.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 E.7: Uživatel ukončil „vysílání“ rozšířených zpráv o svém stavu.

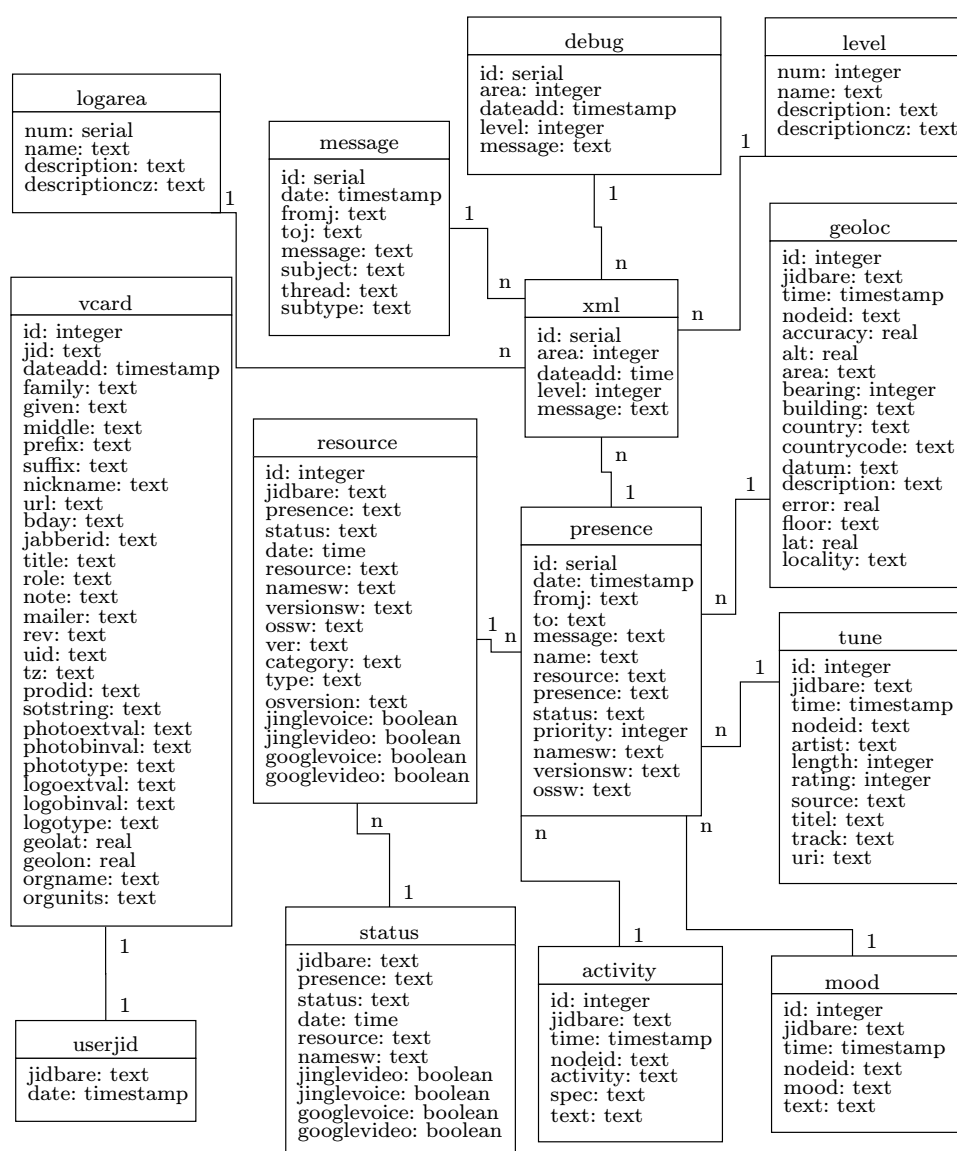
1018 Server informuje všechny účastníky odběru zprávou, která má položku *tune* prázdnou.
1019 Tak jak to prezentuje příklad E.8.
1020

```
1 <message from="user@jabbim.com"
2   to="jabinfo@jabbim.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 E.8: Server informuje klienty o ukončení šíření rozšířeného statusu uživatele *user@jabbim.com*.

1021 Příloha F

1022 Návrh databáze



Obrázek F.1: Struktura databáze

1023 Příloha G

1024 Přehled klientů a jejich rozšíření

Klient	OS	XEP--60	XEP--163	XEP--80	XEP--92	XEP--107	XEP--108	XEP--118	XEP--	XEP--
Adium										
Agile Messenger										
AQQ										
Aytm										
beejive										
Beem										
BitlBee										
Bombus										
BuddyMob										
Chatopus										
Citron										
Claros Chat										
climm										
Coccinella										
Crosstalk										
Digsby										
eM Client										
emite										
Empathy										
Exodus										
Finch										
Gajim										
Galaxium										
glu										
GNU Freetalk										
Gossip										
iChat										
iJab										
IM+										
imov Messenger										
irssi-xmpp										
Jabbear										
Jabber Mix Client										
jabber.el										
Jabbim										
Jabbim for Android										
Jabiru										
JAJC										
Jappix										

Klient	OS	XEP--60	XEP--163	XEP--80	XEP--92	XEP--107	XEP--108	XEP--118	XEP--	XEP--
JBuddy Messenger										
Jeti										
Jitsi (SIP Communicator)										
JWChat										
Kadu										
Kopete										
Lampiro										
m-im										
mcabber										
mChat										
Miranda IM										
Monal IM										
OctroTalk										
OneTeam										
OneTeam for iPhone										
Oyo										
Pandion										
Poezio										
Pidgin										
Prodromus										
Psi										
Psi+										
Quiet Internet Pager (QIP)										
qutIM										
saje										
SamePlace										
Sim-IM										
Slimster										
SoapBox Communicator										
Spark										
SparkWeb										
Synapse										
Talkonaut										
Tigase Messenger										
Tigase Minichat										
Tkabber										
Tlen										
Trillian										
TrophyIM										
V&V Messenger										
Vacuum-IM										
Vayusphere										
WTW										
Xabber										
xmppchat										
Yambi										
Yaxim										

Tabulka G.1: Přehled podporovaných rozšíření u jednotlivých klientů.