

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 Jabberového 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

22. dubna 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	3
2.1	Architektura	3
2.2	XML	6
2.3	Stanza	6
2.4	Rozšíření	9
3	Data mining	12
3.1	Metody dolování dat	14
3.2	Shlukování	16
3.3	Programy	19
4	Implementace	20
4.1	Databáze	20
4.2	Architektura	20
4.3	Robot	20
5	Vyhodnocení výsledků	22
6	Závěr	23
A	Obsah CD	26
B	Manual	27
C	Konfigurační soubor	28
D	Slovník výrazů	29
E	Stanza - základní schéma	30
E.1	Iq	30
E.2	Message	30
E.3	Presence	31
E.4	Přehled průběhu rozšíření	32
F	Přehled klientů a jejich rozšíření	35

¹ Kapitola 1

² Úvod

³ Dnes mezi velmi se rozšiřující technologie na poli síťo Cílem této práce je získat neznámé
⁴ informace z real-time komunikační sítě Jabber.

5 Kapitola 2

6 XMPP

7 V následující kapitole jsou, pro usnadnění a jednodušší pochopení, rozebrány základní sta-
8 vební kameny protokolu Extensible Messaging and Presence Protocol (XMPP). Konkrétně
9 jsou zde popsány stávající vlastnosti implementace, architektura protokolu XMPP obecně
10 [22, 23] a další detaily protokolu [1, 24, 14]. Vzhledem k požadavkům na dolování v da-
11 tech popsaných v následující kapitole je kladen důraz na vybraná rozšíření [21, 12]. Tato
12 rozšíření tvoří základ pro některé rozšířené statusy¹. Další informace použité pro popis a
13 pochopení XML jazyka byly čerpány z [10, 9].

14 Vznik samotného protokolu XMPP je datován do roku 2004 (březen), kdy na něj byl
15 přejmenován Jabber. Původní projekt Jabber byl vytvořen roku 1998 autorem Jeremie
16 Millerem, který ho založil za účelem vytvořit svobodnou otevřenou IM službu. Uvedený
17 projekt měl obsahovat tři základní vlastnosti, do kterých se zahrnují jednoduchost a sro-
18 zumitelnost pro implementaci, jednoduchost v oblasti šíření a otevřenost podobě veřejně
19 dostupného popisu samotného protokolu. Základní vlastnosti a výhody klientů a serverů
20 budou podrobněji popsány níže. Roku 1999, 4.ledna byl vytvořen první server se jménem
21 Jabber. Komunita vývojářů se chopila iniciativy a vytvořila klienty, kteří dokázali se ser-
22 verem komunikovat, pro různé platformy (Linux, Macintosh, Windows). Roku 2004 byl
23 protokol XMPP přidán mezi RFC² dokumenty. Základní norma popisující obecnou struk-
24 turu protokolu je RFC 3920 [22] a RFC 3921 [23], který se zaměřuje na samotný instant
25 messaging a zobrazení stavu. Další zdokumentovaná rozšíření jsou vydávána v podobě tzv.
26 XEP (XMPP Extension Protocol) dokumentů, které jsou známé také pod starším názvem
27 JEP (Jabber Enhancement Proposal). Dnešní počet těchto norem se blíží k číslu 300. Každý
28 XEP obsahuje stav vývoje (schválení), ve kterém se zrovna nachází.

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

32 2.1 Architektura

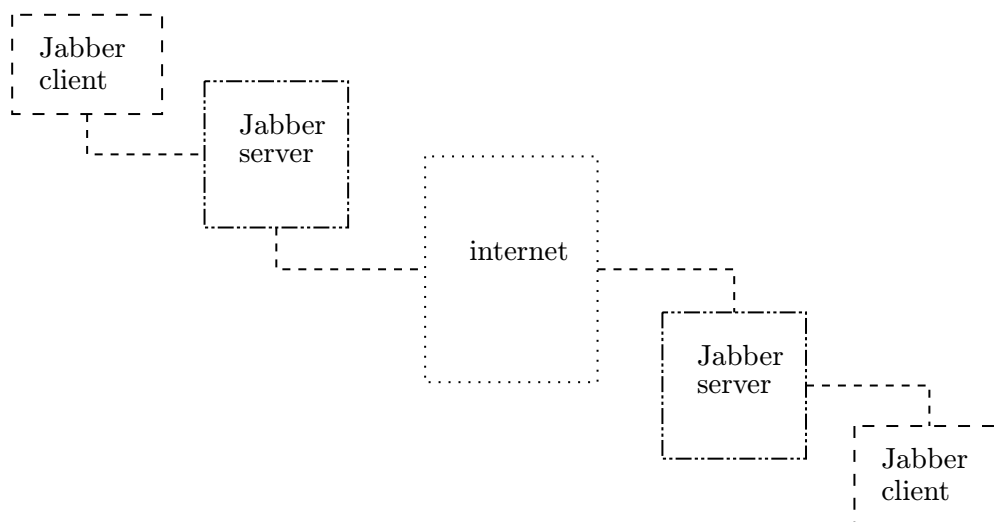
33 Dobře navržená architektura tvoří základ pro správně fungující internetovou technologii.
34 Pro její popis v této práci bylo čerpáno z [1, 24]. Tato struktura se nejvíce podobá struktuře
35 posílání e-mailů. Hlavní předností Jabber sítě je, tak jako u elektronické pošty, její decen-
36 tralizace. V případě Jabberu je decentralizace chápána jako možnost provozovat vlastní

¹User Tune, User Mood, User Location a další

²RFC request of comments – žádost o komentáře

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]. 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.

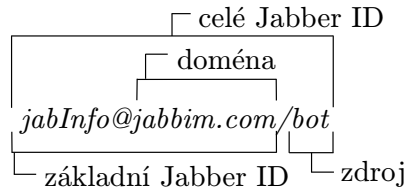
Architektura Jabber serverů využívá velké množství mezi-doménových připojení podobně jako internetový systém elektronické pošty. Komunikace klienta z jedné domény s klientem z jiné na rozdíl od e-mailového modelu nevyžaduje spolupráci třetích stran. Klient se spojí s „domácím“ serverem, který přímo naváže spojení se serverem požadovaného klienta. Tyto vlastnosti jsou zárukou pro bezpečný přenos zpráv, znemožňující „krádeže“ JID³, 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* nebo-li 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.

Další vymoženost JabberID oproti e-mailové adrese je jeho možnost používat prakticky

³uživatelské jméno



Obrázek 2.2: Rozebraná struktura Jabber ID.

libovolné národní znaky u doménových jmen a uživatelských účtů [24]. Využíváním kódování UNICODE, se XMPP stává plně mezinárodní a není jako jiné protokoly omezen rozsahem ASCII tabulky. Přestože je tato vymoženost k dispozici, doposud není žádným výrazným způsobem využívána.

Klient

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

1. komunikace se serverem pomocí TCP soketu
2. rozparsování a následná interpretace příchozí XML zprávy „stanza“ (kapitola 2.3)
3. porozumění sadě zpráv z Jabber jádra

Server

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

Server Jabber je považován za „deamona“, který spravuje tok dat mezi jednotlivými komponentami, které společně tvoří Jabber služby. Například *Jabber Session Manager* (JSM) poskytne funkce pro IM komunikaci a práci se seznamem kontaktů. Komunikace mezi jednotlivými servery je zprostředkována za pomoci komponenty *S2S* (server to server). Při připojení klienta k serveru je komunikace řízená pomocí *C2S* (client to server). Jak již bylo řečeno Jabber síť využívá doménová jména místo špatně zapamatovatelných IP adres. Pro tento způsob identifikace je určena služba *dnsrv*, která názvy překládá.

V tabulce 2.1 jsou shrnuty informace o serverech Jabberu. První sloupec tvoří jméno, následuje programovací jazyk v němž je napsán. Většina aplikací pro servery je vydávána pod licencí GPL⁴. U všech aplikací byla zkoumána nejaktuálnější verze. Její číslo lze nalézt ve třetím sloupci. Všechny servery lze provozovat na operačním systému Linux a Windows.

⁴General Public License-všeobecná veřejná licence GNU

Na platformě Mac OS mohou být použity všechny zde jmenované vyjma jabberd2. Pět z šesti zde představených programů pro server Jabber jsou ve stále vyvíjeny, tedy kromě jabberd14. Hlavním účelem tabulky je prezentovat důležité vlastnosti serverů v oblasti podpory rozšířených statusů. Jedná se o standardy *pubsub*⁵ (XEP-0060) [12] a o jeho verzi zaměřenější více na uživatele *pep*⁶ (XEP-0163) [21]. Obě tato rozšíření tvoří nezbytnou základnu pro *rozšířené statusy* a proto je jejich podpora jak u serverů tak klientů vyžadována. Podrobněji toto téma bude rozebráno v některé následující podkapitole.

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

Tabulka 2.1: Přehled Jabber serverů.

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

2.2 XML

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

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

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, nebo-li 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

⁵Publish-Subscribe

⁶Personal 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.

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 označována pomocí dvou počátečních písmen *Info/Query* nebo-li *IQ*. Na rozdíl od elementu *message* tvoří *iq* spolehlivější přenos, optimalizovaný pro výměnu dat (binární data). K dalším rozdílům patří povinnost příjemce odpovědět na každou přijatou zprávu, nebo-li potvrdit její doručení. Skutečnost, že je na právě požadovanou zprávu odpovězeno, zajišťuje parametr *id*. *Iq* dotaz nebo odpověď musí obsahovat stejnou hodnotu tohoto atributu jako

153 zpráva vytvořená žádajícím subjektem. Další povinný atribut rozděluje iq na čtyři typy.
154 Jednotlivé žádosti na proces nebo akci jsou posílány samostatně [23]. V příloze E je uve-
155 dena rozsáhlejší struktura tohoto elementu. Použití nachází v případech, které nastavují,
156 žádají nebo informace posílají. Tato struktura je využívána pro novou registraci, posílání
157 kontakt listu a další.

158 Příklad 2.3 znázorňuje základní použití elementu iq. Uživatel user posílá dotaz na získání
159 kontakt listu (řádek 5.).

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

Příklad 2.3: Použití elementu iq.

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

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

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

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

Příklad 2.4: Použití elementu *presence*.

177 Obsáhlejší struktura elementu *presence* je zobrazena v příloze E, kde je rovněž k nalezení
178 přehled všech možných stavů.
179

180 Jak již bylo zmíněno v části o Jabber ID Jabber podporuje práci s více současně připoje-
181 nými klienty k jednomu Jabber účtu. Vysvětlení funkčnosti bude prezentováno na příkladu
182 uživatele přihlášeného na stolním počítači a z klienta v mobilním telefonu. U obou těchto
183 připojení je použit stejný Jabber bare, ale odlišného resource, například *domov* a *mobile*.
184 Právě tento rozdíl v tzv. „full“ adrese účtu zajišťuje jednu ze dvou možných podmínek

185 pro správnou adresaci zpráv. Druhá možnost, která bude uplatněna při použití adresy účty
186 pouze ve formě Jabber bare, je nastavení priority u jednotlivých programů. Priorita je číslo
187 v rozsahu hodnot od -128 do 127, kde klient s větší prioritou má přednost před klientem s
188 nižší. Nastane-li případ připojení více klientů se stejnou prioritou, každý server se při roze-
189 sílání zpráv zachová podle vlastní implementace. Některé rozešlou zprávy všem klientům,
190 jiné naopak jen poslednímu přihlášenému.

191 2.4 Rozšíření

192 Dále se tato práce zabývá rozšířeními protokolu XMPP o další vlastnosti k jejichž popisu
193 slouží XEP. Pro tuto práci jsou nepostradatelné „statusy“, pro které tvoří základ standardy
194 XEP-0060 [12] a XEP-0163 [21] zkráceně PEP⁷. Obě tato rozšíření umožňují strukturo-
195 vaně pracovat, používat a přenášet další XEP protokoly. Jako příklady relevantní k práci
196 jsou zde uvedeny protokoly *User Location* (kde se uživatel právě nachází) [6], *User Tune*
197 (co uživatel poslouchá za hudbu) [18], *User Mood* () [20] a *User Activity* () [11]. Jsou to
198 tedy protokoly založené na PEP, které vyžadují podporu nejen v klientech, ale i na straně
199 serveru (zobrazuje tabulka 2.1). S touto informací úzce souvisí další protokol XEP-0115
200 [7], který umožňuje zjistit podporované schopnosti klienta, případně které informace je
201 ochoten přijímat. Tato vlastnost bude popsána níže v části zabývající se podporovanými
202 vlastnostmi.

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

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

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

213 V případě úspěšného přijetí *iq* zprávy serverem, každý, kdo se zaregistroval k odebrání
214 rozšířených statusů, obdrží oznámení ve formě *message*. Oznámení bude také doručeno všem
215 resources. Celá zpráva i všechny další náležitosti jsou uvedeny v příloze E.4.

⁷Personal Eventing via Pubsub

216 Podporované vlastnosti

217 Jednotlivá rozšíření protokolu XMPP jsou nepovinná, a proto nemusí být ve všech klient-
218 ských aplikacích podporována. Pro zjištění podporovaných rozšíření se používá XEP-0115
219 Entity Capabilities [7]. Toto rozšíření výrazně snižuje počet a velikost komunikací a přenosů
220 zpráv mezi uživateli. Dotazem zobrazeným na příkladu 2.6 je zjištěna schopnost jednotli-
221 vých klientů, kterou následně server využije pro správné směřování rozšířených statusů.
222 Všechny zde zmiňované rozšíření a protokoly z této kapitoly je možné u každého klienta
223 (seznam klientů obsahuje tabulka v příloze E.4) vyčíst z atributu *ver* (druhá část u atributu
224 *node*), který je vypočítán ze všech podporovaných protokolů klienta, viz [7].

```
1<iq from="user@jabber.com" id="disco1"
2   to="jabinfo@jabber.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.

225 Další rozšíření

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

228 Prvním rozšířením, nad rámec základních vlastností Jabberu, které zde bude podrobněji
229 rozebráno je elektronická verze klasické vizitky nebo-li *VCard*. Jeho specifikací se zabývají
230 dva standardy. Jelikož novější verze XEP dokumentu [13] se v době psaní této práce na-
231 cházelo ve stavu „experimental“, bylo použito verze starší [16]. Jednoduše řečeno je *VCard*
232 struktura, která nese informace o uživateli jako je jméno, příjmení, e-mail, adresa bydliště
233 i zaměstnání a další údaje. Data jsou dále zveřejňována na síti, z čehož vyplývá, že jsou
234 dostupná ostatním uživatelům. Vyplnění těchto osobních údajů je dobrovolné a tak se u
235 některých uživatelů nachází pouze přezdívka a JID, které jsou často předdefinovány auto-
236 maticky. Nedílnou součástí všech sociálních a komunikačních systému je tzv. „Avatar“⁸. V
237 síti Jabber tomu není jinak, a proto je samotný obrázek zahrnut přímo do *VCard* v položce
238 *photo*. Podrobnější informace o jeho nastavení a přijímání je možné nalézt v *vCard-Based*
239 *Avatars* [19], který jej definuje.

240 Díky základní podmínce XMPP protokolu (otevřenost) existuje mnoho různých aplikací
241 pomocí, kterých lze v síti Jabber komunikovat. S programy, používanými uživateli, úzce
242 souvisí další zde implementované rozšíření. Jedná se o realizaci *Software Version* dokumentu
243 [17], který se právě zabývá získáváním informací o samotných aplikacích. Je-li toto rozšíření
244 podporováno je díky němu možné zjistit jméno a verzi používané aplikace. Informace o
245 operačního systému často nejsou kvůli bezpečnosti ani vyplněny. Podrobnější informace o
246 softwarové výbavě klienta je možné zjistit pomocí XEP [7], o kterém již bylo dříve psáno v
247 odstavci zabývajícím se podporovanými vlastnostmi klientských aplikací.

248 S rozšířením tzv. „chytrých“ mobilních zařízení mezi širší veřejnost vzniklo několik no-
249 vých disciplín spojených s určováním zeměpisné polohy jako je například geocaching. Geo-
250 grafická poloha [6] je přenášena jak ve formě GPS⁹ souřadnic tak i ve tvaru „civilním“, jako

⁸fotografie, logo nebo ikona

⁹Global Positioning System

251 je stát, ulice, město nebo číslo poschodí v budově a další. Mnoho aplikací, které mají k dis-
252 pozici GPS přijímač, vysílají a aktualizují zeměpisné informace automaticky, například po
253 určité době nebo změně polohy o určitou vzdálenost. Toto a další níže popsané rozšíření jsou
254 postaveny na již zmiňovaném PEP. Některé části protokolů jsou zjednodušeny a připraveny
255 tím pro „mobilní instant messaging“.

256 Pro sdělení informací o stavu klienta není v základní verzi Jabberu mnoho. Pomocí
257 presence je možné „pouze“ prozradit zda je uživatel připraven komunikovat nebo je mo-
258 mentálně nedostupný a to v několika verzích lišících se délkou nepřítomnosti. Pokročilejší
259 nastavení statusu nabízí *User Mood* [20] a to ve formě sdělení současné nálady jako je
260 například radost. Další možné upřesnění činnosti uživatele jsou definovány v *User Activity*
261 [11], kde každá činnost je složena z povinné obecné kategorie a nepovinné, která informaci
262 upřesňuje. Příkladem může být *eating* a *having-a-snack*.

263 K poslednímu rozšíření implementovanému v této práci patří *User Tune* [18], které
264 umožňuje uživateli šířit informace o aktuálně poslouchané hudbě. Některé dnešních hudební
265 přehrávače dokáží automaticky spolupracovat s IM klientem a předávat informace o hudbě
266 bez nutného lidského zásahu.

267 Výše popsáná rozšíření nejsou z velké míry podporována aplikacemi, což způsobuje ne-
268 spokojenost mnoha uživatelů využívající hlavně mobilní zařízení. Například v předcházející
269 zmiňované části o poslouchané hudbě, při nepodporovanosti programu, je posíláno pomocí
270 normální presence. Jméno skladatele, alba a další podrobnosti jsou shrnuty do statusu,
271 tudíž jsou doručeny všem uživatelům z kontakt listu.

272 Kapitola 3

273 Data mining

274 Třetí kapitola se zabývá procesem dobývání znalostí z databází. Popisuje jej jako disciplínu,
275 která vznikla za účelem vytěžení informací z dat, která jsou v nepřehledném množství uklá-
276 dána v databázích. Díky velikosti dnešních disků, objem ukládaných dat neustále roste. S
277 tím také úzce souvisí zvětšující se poměr nepotřebných a zašumělých dat vůči užitečným
278 informacím.

279 Na začátku kapitoly je rozebrán pojem získávání znalostí databází, jehož jednu podstat-
280 nou část tvoří samotný data mining. Dále je vysvětlena základní terminologie, pro kterou bylo
281 čerpáno z [8]. Celá první podkapitola je věnována vybraným metodám pro dolování dat a
282 vlastnostem, které je od sebe navzájem odlišují. Jsou zde rozebrány *asociační pravidla*, pro
283 jejichž popis bylo čerpáno z [2]. Pro ostatní metody, které jsou popsány dále, byla jako zdroj
284 informací použita kniha [5]. Poté následuje druhá podkapitola, která se podrobněji zabývá
285 jednou z metod pro dolování dat a to *shlukováním*. Obsahem této části jsou již konkrétní
286 algoritmy pro shlukování dat [25, 3] a také metoda *k-Means* využívaná v praktické části
287 této práce. Kapitulu uzavírá přehled vybraných programů pro data mining a podrobnější
288 seznámení s programem *RapidMiner*, který je v této práci využíván pro samotné dolování.

289 Terminologie

290 Pojem data mining nebo-li česky dolování dat se začal ve vědeckých kruzích objevovat
291 počátkem 90. let 20. století. První zmínka pochází z konferencí věnovaných umělé in-
292 teligenci (IJCAI'89¹–mezinárodní konference konaná v Detroitu, AAAI'91² a AAAI'93–
293 americké konference v Californii a Washingtonu, D.C) [2].

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

- 300 1. počet záznamů nebo-li objektů v databázi
- 301 2. počet polí nebo-li atributů objektů v databázi

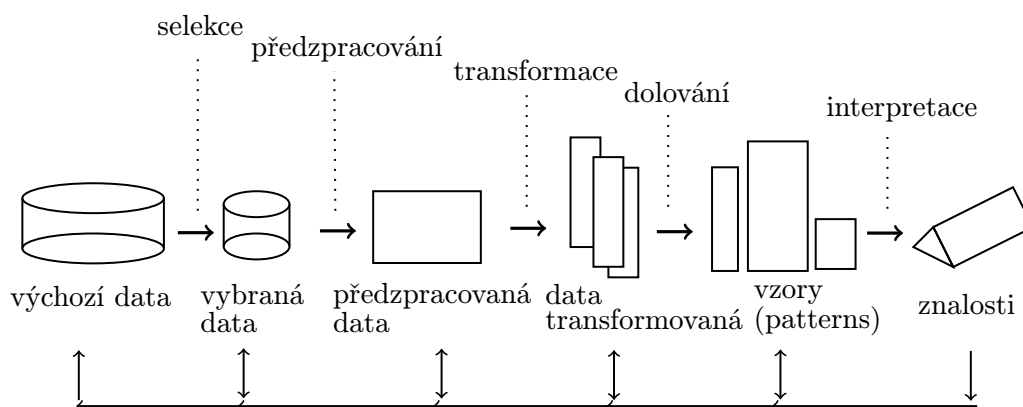
¹International Joint Conference on Artificial Intelligence

²Association for the Advancement of Artificial Intelligence

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

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

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



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

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

- **čištění dat** – fáze, ve které jsou nepodstatné údaje odstraněny z kolekce.
- **integrace dat** – kombinování heterogenních dat z několika zdrojů do společného jediného zdroje.
- **výběr dat** – rozhodování o relevantních datech.
- **transformace dat** – také známý jako konsolidace dat. Fáze, ve které jsou vybraná data transformována do formy vhodné pro dolování.

³Knowledge Discovery in Database

- 328 • **data mining** – zásadní krok, ve kterém jsou aplikovány vzory na data.
- 329 • **hodnocení modelů** – vzory dat zastupují získané znalosti.
- 330 • **prezentace znalostí** – konečná fáze, zjištěné poznatky jsou reprezentovány uživa-
331 teli. Tento základní krok využívá vizualizační techniky, které pomáhají uživa-
332 telům porozumět a správně interpretovat získané výsledky.

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

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

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

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

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

346 Jedna z několika existujících definic pojmu data je uvedena v definici č. 3.0.3, která je
347 popisuje z pohledu informačního. Data často nemají sémantiku (význam) a bývají zpraco-
348 vána čistě formálně.

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

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

355 3.1 Metody dolování dat

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

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

- 364 • schopnosti reprezentace shluků (např. otázka lineární separability)
- 365 • srozumitelnosti nalezených znalostí pro uživatele (symbolické vs. subsymbolické me-
366 tody)

- efektivnosti znovupoužití nalezených znalostí
- vhodnosti typů dat
- a další ...

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

Asociační pravidla

Při popisu asociačních pravidel, která jsou založena na syntaxi *IF-THEN*, bylo čerpáno z [2]. Jejich rozšíření se datuje do 90. let 20. století, kdy byly panem Agrawalem představeny v souvislosti s analýzou „nákupního košíku“.

Použitelnost bude vysvětlena právě na příkladu analýzy nákupního košíku. Podstata příkladu je tvořena zákazníkem a jeho systémem nakupování. Jsou zjišťovány produkty, které jsou nakupovány současně. Hledají se nebo-li jsou vytvářeny společné vazby (asociační pravidla) mezi výrobky a určuje se jejich spolehlivost. Na základě těchto závislostí je upravováno umístění jednotlivých výrobků.

Obecně jsou tedy asociační pravidla považována za konstrukci, která z hodnot jedné transakce odvozuje možnost výskytu závislostí v jiných transakcích. Jsou tedy hledány všechny vnitřní závislosti existující mezi daty.

Podle knihy Berky [2] je základní myšlenka asociačních pravidel *IF-THEN* převedena do jiné terminologie:

$$\text{Ant} \Rightarrow \text{Suc}$$

kde *Ant* bývá interpretován jednou možností z výčtu – předpoklad, *IF*, levá strana pravidla nebo antecedent a *Suc* je chápán jako – závěr, *ELSE*, pravá strana pravidla, sukcedent. Níže jsou uvedeny základní vlastnosti:

$$n(\text{Ant} \wedge \text{Suc}) = \mathbf{a}; n(\text{Ant} \wedge \neg \text{Suc}) = \mathbf{b}; n(\neg \text{Ant} \wedge \text{Suc}) = \mathbf{c}; n(\neg \text{Ant} \wedge \neg \text{Suc}) = \mathbf{d};$$

$$n(\text{Ant}) = a+b = \mathbf{r}; n(\neg \text{Ant}) = c+d = \mathbf{s}; n(\text{Suc}) = a+c = \mathbf{k}; n(\neg \text{Suc}) = b+d = \mathbf{l}; n = a+b+c+d;$$

všechna pravidla jsou shrnuta v tabulce 3.1, z nichž jsou dále počítány různé charakteristiky a následně tak hodnoceny zjištěné znalosti.

	Suc	¬Suc	Σ
Ant	a	b	r
¬Ant	c	d	s
Σ	k	l	n

Tabulka 3.1: Kontingenční tabulka převzata z [2].

Mezi základní charakteristiky asociačních pravidel podle Agrewalova patří *podpora* a *spolehlivost*.

397 Klasifikace

398 Klasifikace bude opět vysvětlena na příkladu, převzatého z [8]. Podle obsahu databáze
399 nebo dotazníku bude každý klient banky zařazen do různých krizových skupin. Na základě
400 těchto skupin pracuje „credit skóring“, jež klientovi poskytne nebo odepře například úvěr v
401 bance. Další příklady využití jsou například ve zdravotnictví. Na základě zdravotního stavu
402 pacienta a jeho příznaků, je pacient zařazen do tříd, které reprezentují jednotlivé nemoci.

403 Klasifikací jsou, podle [5], jednotlivé zkoumané elementy rozděleny (podle hodnot atri-
404 butů) do vhodných kategorií, které jsou předem vytvořeny z navzájem podobných objektů
405 (tvorba profilů třídy). Při této metodě je upřednostňována přesnost před jednoduchostí a
406 rychlostí. Zdroje klasifikovaných objektů jsou většinou tvořeny jednotlivými řádky v data-
407 bázi. Vzory dat vytváří instance, jejichž vlastnosti reprezentují atributy vyjádřené číselnou
408 hodnotou.

409 Modely

410 Základem modelů jsou trénovací data. Níže uvedený příklad vybraných klasifikačních mo-
411 delů, byl čerpán z [5]:

- 412 • Rozhodovací stromy
- 413 • Neuronové sítě
- 414 • Statistické metody
- 415 • Klasifikační pravidla
- 416 • Využití vzdálenosti
- 417 • a další ...

418 Predikce

419 Predikce je řazena mezi velmi známé procesy, které na základě získaných znalostí předpoví-
420 dají následující vývoj. Chronologicky seřazená data a vývoj jejich hodnot v minulosti tvoří
421 základ pro určení hodnot budoucích. Předpokládá se, že na základě informací získaných z
422 dat v minulosti, bude možné postavit modely, které se budou chovat stejně nebo alespoň
423 podobně i v budoucnu. Využití naleznou v předpovědi počasí (z naměřených meteorologic-
424 kých hodnot se určují budoucí předpokládané teploty), při vývoji cen na burze a dalších.
425 Podklady pro popis predikce byly čerpány z [2, 5].

426 Shlukování

427 Metoda zaměřená na dělení objektů do předem neznámých skupin. Proces dělení probíhá
428 na základě specifikace objektů a jejich odlišnosti od ostatních shluků. Tato část, pro kterou
429 bylo čerpáno z [25, 3], bude podrobně rozebrána v následující podkapitole.

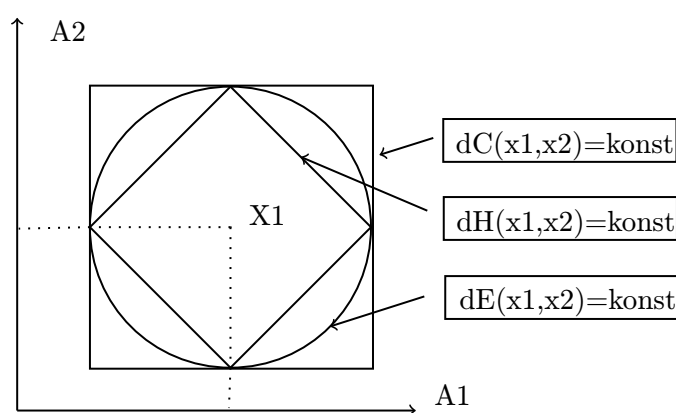
430 3.2 Shlukování

431 V této podkapitole je shlukování rozděleno na několik metod shlukové analýzy podle [5].
432 U každé z nich jsou popsány její základní vlastnosti a uvedeny nejrozšířenější algoritmy.

433 Poslední metoda *metoda rozkladu* je rozebrána podrobněji z důvodu jejího praktického
 434 využití v této práci.

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

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



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

447 Metody založené na modelu

448 Metody založené na modelu se pokouší přiřadit data k určitému matematickému modelu na
 449 základě společných optimalizovaných vlastností. Většina procesů je založena na předpokladu
 450 generování dat pomocí standardních statistik.

451 Mezi zástupné metody této shlukovací analýzy se řadí Expectation–Maximization (EM)
 452 a SOON⁴. Algoritmus SOON je založen na neuronové síti. Je to metoda vycházející z
 453 algoritmu SOM⁵ [25]. Metoda EM je rozšířením algoritmu *k-means*, který bude podrobně
 454 rozebrán v následující části.

455 Metody hierarchické

456 Hlavní princip metody hierarchického shlukování je založen na tvorbě stromové hierarchie
 457 shluků, která je známá pod názvem *dendrogram*. Hierarchické metody, podle [5], mohou být
 458 rozděleny do dvou skupin a to na základě principu, kterým jsou dendrogramy vytvářeny.
 459 První možnost je *aglomerativní přístup*, který shlukuje menší shluky, kdy výsledkem je jen

⁴Self Organizing Oscillator Network

⁵Self-Organizing Map

460 jeden. Druhý přístup, *divizní*, je založen na opačném předpokladu. Tedy že na počátku
461 je jeden velký shluk, který je postupně rozdělován dokud není počet shluků roven počtu
462 objektů [25]. Mezi zástupce této metody například patří algoritmus AGNES⁶.

463 Metody založené na mřížce

464 Metody založené na mřížce kvantují datový prostor do konečného počtu pravoúhlých buněk,
465 které jsou uspořádány do víceúrovňové mřížkové struktury. Zmíněná struktura tvoří základ
466 pro shlukové operace. Hlavní výhoda tohoto přístupu je rychlost zpracování, které většinou
467 nebere ohled na počet datových objektů. Čas zpracování závisí pouze na počtu buněk v
468 každé dimenzi kvantovaného prostoru.

469 Mezi zástupce metod založených na mřížce patří metoda STING⁷, který pracuje se
470 statickými informacemi uloženými v buňkách mřížky. Algoritmus je rozdělen do dvou částí.
471 První si klade za cíl rekurzivně rozdělit datový prostor na pravoúhlé buňky. Druhá fáze
472 testuje spojitost mezi sousedy relevantních buněk [25].

473 Mezi další metody založené na mřížce patří WaveCluster⁸, využívající vlnkové transfor-
474 mace k rozdělení prostoru dat. Tato transformace zdůrazňuje shluky v prostoru a objekty
475 jim vzdálené potlačuje [5].

476 Metody založené na hustotě

477 Vychází z m -rozměrného prostoru, ve kterém jsou zobrazeny objekty ve formě bodů. Místa
478 v prostoru s větší koncentrací objektů ve srovnání s ostatními oblastmi jsou nazývány
479 shluky. Výchozí předpoklad je existence okolí jednotlivých bodů (sousedství). Jedna z cha-
480 rakteristik metod založených na hustotě je schopnost vypořádat se s vzdálenými hodnotami,
481 označovanými jako šum [25].

482 Jako příklad je uvedena metoda DBSCAN⁹, která je založena na hustotě objektů v
483 prostoru. U jednotlivých objektů je zkoumáno jejich okolí. Algoritmus je ovlivňován dvěma
484 parametry ε (velikost shluku) a $MinPts$ (minimální počet objektů v daném shluku), které
485 spolu úzce souvisí (viz [5]). Bod splňující obě podmínky je označen za jádro. Za pomoci
486 jader je rozšiřována množina objektů spojených na základě hustoty. Obsahuje-li jádro x_1
487 ve svém okolí další jádro x_2 znamená to, že jádro x_1 je přímo dosažitelné z jádra x_2 . Tímto
488 způsobem jsou vytvářeny výsledné *shluky*. V opačném případě, body, které nesplňují dvě
489 zmíněné podmínky, jsou označeny jako *šum*.

490 Shlukování velkých dat

491 Všechny zde doposud zmiňované metody poskytují dobré výsledky pouze s malým počtem
492 dimenzí, tak jak je to popsáno v [8]. S narůstajícím počtem atributů roste počet nerelevant-
493 ních dimenzí určených pro shlukování. S tímto také přibývá zvětšená produkce zašumění a
494 znesnadnění nalezení relevantních shluků. Data jsou roztroušena do mnoha dimenzí a tím
495 odpadá možnost použití vzdálenostních funkcí.

496 Zmíněné problémy shlukování velkých dat řeší dvě techniky *metoda transformace rysů*
497 *a metoda výběru atributů*. Pro efektivní shlukování je možné použít například algoritmus

⁶AGglomerative Nesting

⁷STatistical INformation Grid

⁸Clustering Using Wavelet Transformation

⁹Density-Based Spatial Clustering of Applications with Noise

498 CLIQUE¹⁰.

499 **Metody rozkladu**

500 ...

501 **k –Means**

502 ...popis a za pis algoritmu ...

503 **3.3 Programy**

504 **Fitminer**

505 **Rapid miner**

506 **Weka**

¹⁰CLustering In QUES

507 Kapitola 4

508 Implementace

509 4.1 Databáze

510 Návrh databáze

511 4.2 Architektura

512 –robot plus databaze –rapidminer–spusteni davkove/vlastni algoritmus v PHP –webova
513 implementace prezentuje vysledky

514 4.3 Robot

515 –class diagram, pomoci nej popsat strukturu, mozne rozsireni

516 gloox

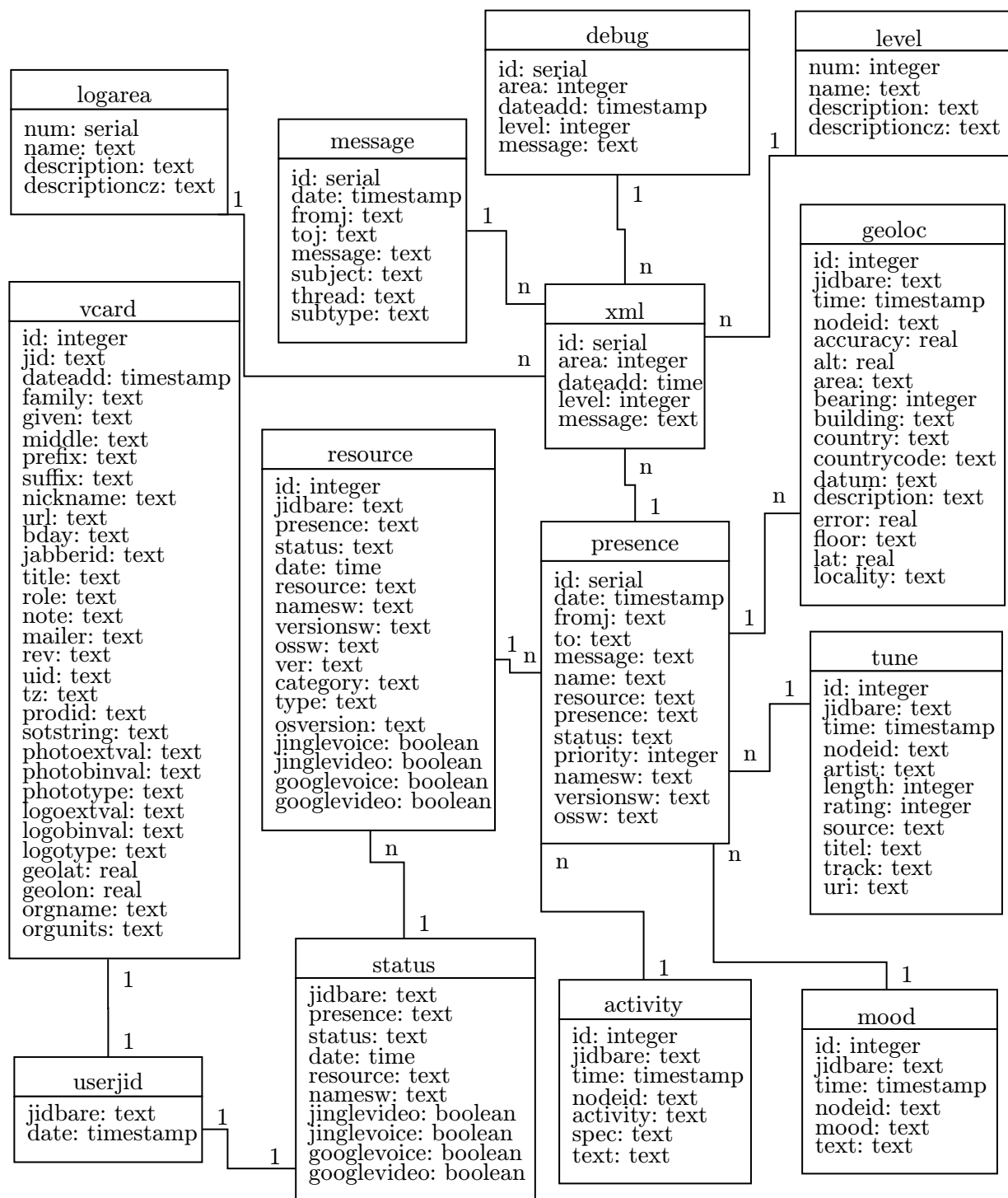
517 Gloox je stabilní Jabber/XMPP knihovna vydávána pod licencí GNU GPL. Je určena pro
518 vývoj klienta a komponent. Jelikož je psána v ANSCI C++ mezi podporované platformy
519 patří Linux, Windows, Mac OS X, Symbian/Nokia S60, FreeBSD a další systémy stavějící
520 na ANSI C++ kompilátor.

521 Pomocí knihovny gloox je psán bot v této práci. Byla vybrána na základě požadavku
522 psaní programu v jazyce C/C++ a operačním systémem Linux. V porovnání s jinými knihov-
523 nami pro jazyk C nebo C++ disponuje lepší podporou a dokumentací. Gloox plně podpo-
524 ruje standart XMPP Core [22] a z větší části i standard XMPP IM [23]. Dodatečně je plně
525 podporováno kolem 30 XEP standardů například XEP-0054 vcard-temp a další.

526 Návrh bota

527 Transformace

528 temporalni databaze



Obrázek 4.1: Struktura databáze

529 Kapitola 5

530 Vyhodnocení výsledků

531 –charakter dat, jaka data jsem nasbiral –

532 Kapitola 6

533 Závěr

534

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.
22. dubna 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. 22. dubna 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.
22. dubna 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.
22. dubna 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. 22. dubna 2011].
URL <http://xmpp.org/extensions/xep-0060.html>

- 565 [13] Mizzi, S.; Saint-Andre, P.: XEP-0292: vCard4 Over XMPP. [online], 02-26-2008, [cit.
566 22. dubna 2011].
567 URL <http://xmpp.org/extensions/xep-0292.html>
- 568 [14] Moore, D.; Wright, W.: *Jabber developer's handbook*. Indianapolis: Sams Publishing,
569 první vydání, 2004, 487 s., iISBN 06-723-2536-5.
- 570 [15] Nemrava, M.; Pospíšil, J.: Dolování dat a jeho aplikace. [online], 2006, [cit. 22. dubna
571 2011].
572 URL http://www.spatial.cs.umn.edu/paper_ps/dmchap.pdf
- 573 [16] Saint-Andre, P.: XEP-0054: vcard-temp. [online], 07-16-2008, [cit. 22. dubna 2011].
574 URL <http://xmpp.org/extensions/xep-0054.html>
- 575 [17] Saint-Andre, P.: XEP-0092: Software Version. [online], 02-15-2007, [cit. 22. dubna
576 2011].
577 URL <http://xmpp.org/extensions/xep-0092.html>
- 578 [18] Saint-Andre, P.: XEP-0118: User Tune. [online], 30-01-2008, [cit. 22. dubna 2011].
579 URL <http://xmpp.org/extensions/xep-0118.html>
- 580 [19] Saint-Andre, P.: XEP-0153: vCard-Based Avatars. [online], 16-08-2006, [cit.
581 22. dubna 2011].
582 URL <http://xmpp.org/extensions/xep-0153.html>
- 583 [20] Saint-Andre, P.; Meijer, R.: XEP-0107: User Mood. [online], 29-10-2008, [cit.
584 22. dubna 2011].
585 URL <http://xmpp.org/extensions/xep-0107.html>
- 586 [21] Saint-Andre, P.; Smith, K.: XEP-0163: Personal Eventing Protocol. [online],
587 12-07-2010, [cit. 22. dubna 2011].
588 URL <http://xmpp.org/extensions/xep-0163.html>
- 589 [22] Saint-André, P.: Extensible Messaging and Presence Protocol (XMPP): Core.
590 [online], 10-2004, [cit. 22. dubna 2011].
591 URL <http://tools.ietf.org/html/rfc3920>
- 592 [23] Saint-André, P.: Extensible Messaging and Presence Protocol (XMPP): Instant
593 Messaging and Presence. [online], 10-2004, [cit. 22. dubna 2011].
594 URL <http://tools.ietf.org/html/rfc3921>
- 595 [24] Saint-André, P.; Smith, K.; Troncon, R.: *XMPP : the definitive guide : building
596 real-time applications with jabber technologies*. Sebastopol: O'Reilly, první vydání,
597 2009, 287 s., iISBN 978-059-6521-264.
- 598 [25] Řezánková, H.; Húsek, D.; Snášel, V.: *Shluková analýza dat*. Praha: Professional
599 Publishing, druhé vydání, 2009, 218 s., iISBN 978-808-6946-818.

⁶⁰⁰ **Příloha A**

⁶⁰¹ **Obsah CD**

⁶⁰² **Příloha B**

⁶⁰³ **Manual**

⁶⁰⁴ **Příloha C**

⁶⁰⁵ **Konfigurační soubor**

606 Příloha D

607 Slovník výrazů

608 **DNS** — Domain Name System

609 **GPG** — dkshckdsjvlsdjvodsvjdfokj

610 **IM služby** — dkshckdsjvlsdjvodsvjdfokj

611 **IP** — Internet Protocol

612 **JEP** — dkshckdsjvlsdjvodsvjdfokj

613 **JID** — dkshckdsjvlsdjvodsvjdfokj

614 **SASL** — dkshckdsjvlsdjvodsvjdfokj

615 **TCP** — dkshckdsjvlsdjvodsvjdfokj

616 **TLS** — dkshckdsjvlsdjvodsvjdfokj

617 **WWW** — dkshckdsjvlsdjvodsvjdfokj

618 **XEP** — dkshckdsjvlsdjvodsvjdfokj

619 **XML** — dkshckdsjvlsdjvodsvjdfokj

620 **XMPP** — dkshckdsjvlsdjvodsvjdfokj

621 **e-mail** — dkshckdsjvlsdjvodsvjdfokj

622 **jabber** — dkshckdsjvlsdjvodsvjdfokj

623 **klient** — dkshckdsjvlsdjvodsvjdfokj

624 **presence** — dkshckdsjvlsdjvodsvjdfokj

625 **server** — dkshckdsjvlsdjvodsvjdfokj

626 **stanza** — dkshckdsjvlsdjvodsvjdfokj

627 **vCard** — dkshckdsjvlsdjvodsvjdfokj

628 Příloha E

629 Stanza - základní schéma

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

634 E.1 Iq

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

Příklad E.1: Popis elementu *iq*.

635 E.2 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>
```

Příklad E.2: Popis elementu *message*.

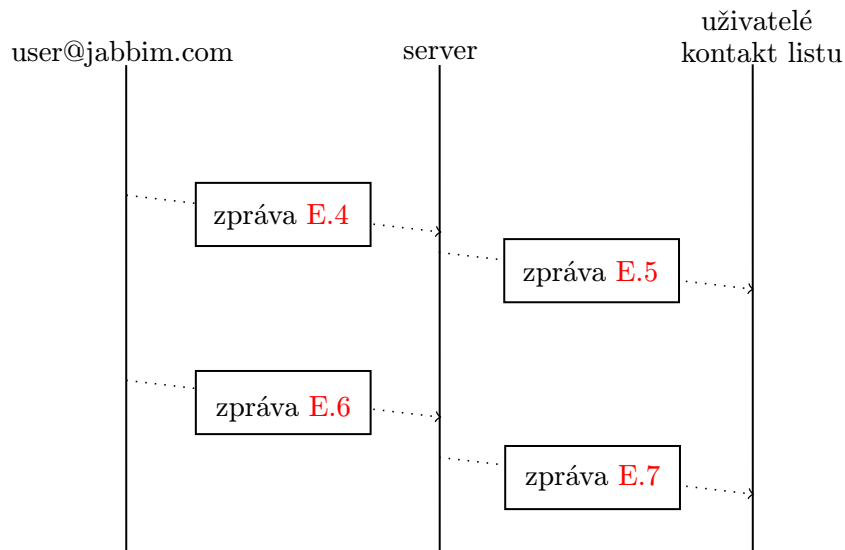
E.3 Presence

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

Příklad E.3: Popis elementu *presence*.

637 E.4 Přehled průběhu rozšíření

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



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

639

```

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>

```

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

640 Server obdrží informace od klienta *user* zprávu o přehrávací hudbě. Pomocí elementu
641 *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>
```

Příklad E.5: Server informuje uživatele podporující rozšíření o stavu *user@jabbim.com*.

642 Zpráva o přehrávané hudbě je také přeposlána všem otevřeným spojením uživatele *user*,
643 ukázáno na příkladě E.6.
644

```
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>
```

Příklad E.6: Server přepoše informace o přehrávané hudbě všem otevřeným spojením
uživatele *user@jabbim.com*.

645 Přestane-li uživatel *user* poslouchat/vysílat informace o přehrávané hudbě, provede to
646 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>
```

Příklad E.7: Uživatel ukončil „vysílání“ rozšířených zpráv o svém stavu.

647 Server informuje všechny účastníky odběru zprávou, která má položku *tune* prázdnou.
648 Tak jak to prezentuje příklad E.8.
649

```
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>
```

Příklad E.8: Server informuje klienty o ukončení šíření rozšířeného statusu uživatele *user@jabbim.com*.

650 Příloha F

651 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 F.1: Přehled podporovaných rozšíření u jednotlivých klientů.