

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
26. 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 | 7 |
| 2.4 | Rozšíření | 9 |
| 3 | Data mining | 12 |
| 3.1 | Metody dolování dat | 14 |
| 3.2 | Shlukování | 16 |
| 3.3 | Programy | 20 |
| 4 | Implementace | 22 |
| 4.1 | Databáze | 22 |
| 4.2 | Architektura | 23 |
| 4.3 | Robot | 23 |
| 5 | Vyhodnocení výsledků | 25 |
| 6 | Závěr | 26 |
| A | Obsah CD | 29 |
| B | Manual | 30 |
| C | Konfigurační soubor | 31 |
| D | Slovník zkratk | 32 |
| E | Stanza - základní schéma | 33 |
| E.1 | Iq | 33 |
| E.2 | Message | 33 |
| E.3 | Presence | 34 |
| E.4 | Přehled průběhu rozšíření | 35 |
| F | Přehled klientů a jejich rozšíření | 38 |

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 statusy, jako je například User Tune [18], User
13 Mood [20], User Location [6] a další. Další informace použité pro popis a pochopení XML
14 jazyka byly čerpány z [10, 9].

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

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

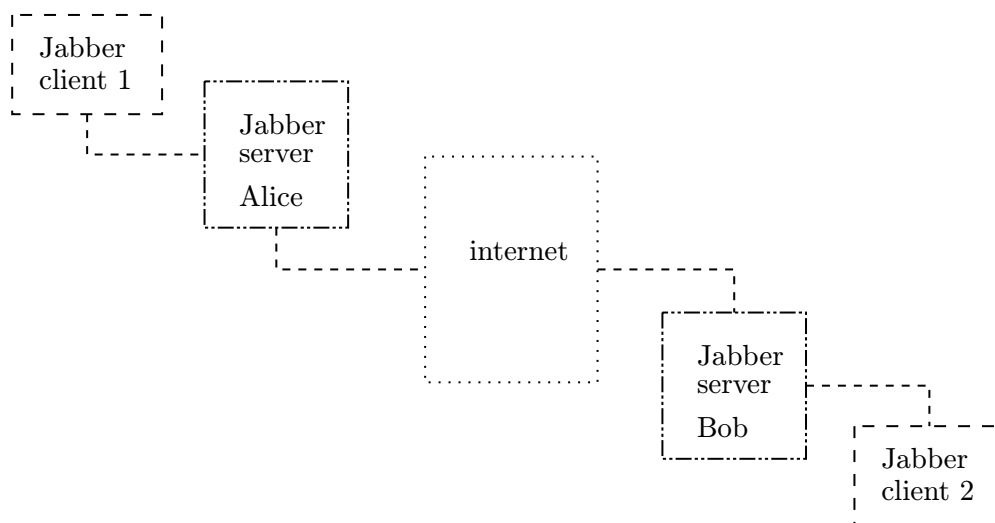
33 2.1 Architektura

34 Dobře navržená internetová technologie je tvořena správně fungujícími komponenty, které mezi
35 sebou dokáží vytvořit spojení a následně započít komunikaci. Pro popis Jabber architek-
36 tury v této práci bylo čerpáno z [1, 24]. Tato struktura se nejvíce podobá struktuře posílání
37 e-mailů. Hlavní předností Jabber sítě je, tak jako u elektronické pošty, její decentralizace.

¹RFC request of comments – žádost o komentáře

38 V případě Jabberu je decentralizace chápána jako možnost provozovat vlastní server, na
 39 rozdíl od jiných komunikačních systémů jako je například facebook, kde existuje pouze je-
 40 diný poskytovatel služby. V případě serveru je kladen důraz na spolehlivost a rozšiřitelnost
 41 a u klienta na uživatele. Každý server pracuje samostatně, což znamená, že chod ani vý-
 42 padek jiné datové stanice žádným způsobem jeho běh neovlivní. V případě výpadku jiného
 43 serveru bude nedostupný pouze seznam kontaktů a služeb, které registrovaným uživatelům
 44 poskytoval.

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



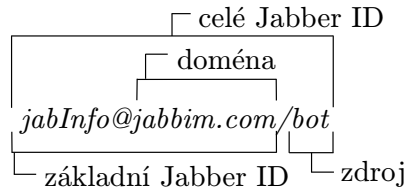
Obrázek 2.1: Distribuovaná architektura Jabber.

47 Architektura Jabber serverů využívá velké množství mezi-doménových připojení po-
 48 dobně jako internetový systém elektronické pošty. Komunikace klienta z jedné domény s
 49 klientem z jiné na rozdíl od e-mailového modelu nevyžaduje spolupráci třetích stran. Kli-
 50 ent se spojí s „domácím“ serverem, který přímo naváže spojení se serverem požadovaného
 51 klienta. Tyto vlastnosti jsou zárukou pro bezpečný přenos zpráv, znemožňující „krádeže“
 52 JID², který je popsán níže, a spamování.

54 Jabber ID

55 Jabber ID (JID) je jednoznačný virtuální identifikátor uživatele na síti. V případě založení
 56 účtu nejsou rozlišována velká a malá písmena, což znamená, že Jabber není case-sensitive.
 57 Jednoznačný Jabber identifikátor je složen ze dvou částí *Jabber bare* nebo-li čisté ID a
 58 *resource* [22]. Základní část na první pohled připomíná e-mailovou adresu *user@server*.
 59 Druhá část slouží k přesné identifikaci jednotlivých spojení. Je použita ke směrování síťového
 60 provozu s uživateli v případě otevření většího množství spojení pod jedním uživatelem.
 61 Společně Jabber bare a resource tvoří tzv. *full JID* — *user@server/resource* například
 62 *jabInfo@jabim.cz/bot*. Jednotlivé části uživatelského jména popsané v tomto odstavci jsou
 63 ukázány v obrázku 2.2.

²uživatelské jméno



Obrázek 2.2: Rozebraná struktura Jabber ID.

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

69 Klient

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

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

80 Server

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

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

V tabulce 2.1 jsou shrnuty informace o serverech Jabberu. První sloupec tvoří jméno, následuje programovací jazyk v němž je napsán. Většina aplikací pro servery je vydávána pod licencí GPL³. U všech aplikací byla zkoumána nejaktuálnější verze. Její číslo lze nalézt ve třetím sloupci. Všechny servery lze provozovat na operačním systému Linux a Windows. Na platformě Mac OS mohou být použity všechny zde jmenované vyjma jabberd2. Pět z šesti zde představených programů pro server Jabber jsou ve stále vyvíjeny, tedy kromě jabberd14. Hlavním účelem tabulky je prezentovat důležité vlastnosti serverů v oblasti podpory rozšířených statusů. Jedná se o standardy *pubsub*⁴ (XEP-0060) [12] a o jeho verzi zaměřenější více na uživatele *pep*⁵ (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 |
| 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.

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

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.

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, 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 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@jabbbim.com"
2              to="jabinfo@jabbbim.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 zpráva vytvořená žádajícím subjektem. Další povinný atribut rozděluje iq na čtyři typy. Jednotlivé žádosti na proces nebo akci jsou posílány samostatně [23]. V příloze E je uvedena rozsáhlejší struktura tohoto elementu. Použití nachází v případech, které nastavují, žádají nebo informace posílají. Tato struktura je využívána pro novou registraci, posílání seznamu kontaktů a další.

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

```

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

```

Příklad 2.3: Použití elementu *iq*.

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

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

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

```

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

```

Příklad 2.4: Použití elementu *presence*.

Obsáhlejší struktura elementu *presence* je zobrazena v příloze E, kde je rovněž k nalezení

185 přehled všech možných stavů.

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

197 2.4 Rozšíření

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

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

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

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

223 Podporované vlastnosti

224 Jednotlivá rozšíření protokolu XMPP jsou nepovinná, a proto nemusí být ve všech klient-
225 ských aplikacích podporována. Pro zjištění podporovaných rozšíření se používá XEP-0115
226 Entity Capabilities [7]. Toto rozšíření výrazně snižuje počet a velikost komunikací a přenosů

⁶Personal Eventing via Pubsub

```

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

```

Příklad 2.5: Začátku vysílání rozšířeného statusu.

227 zpráv mezi uživateli. Dotazem zobrazeným na příkladu 2.6 je zjištěna schopnost jednotlivých klientů, kterou následně server využije pro správné směrování rozšířených statusů. 228
 229 Všechny zde zmiňované rozšíření a protokoly z této kapitoly je možné u každého klienta (seznam klientů obsahuje tabulka v příloze E.4) vyčíst z atributu *ver* (druhá část u atributu 230 *node*), který je vypočítán ze všech podporovaných protokolů klienta, viz [7]. 231

```

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.

232 Další rozšíření

233 V následujících několika odstavcích budou přiblíženy specifikace jednotlivých rozšíření XEP, 234 které slouží jako zdrojová data pro dolování a jsou relevantní k tématu práce.

235 Prvním rozšířením, nad rámec základních vlastností Jabberu, které zde bude podrobněji 236 rozebráno je elektronická verze klasické vizitky nebo-li *VCard*. Jeho specifikací se zabývají 237 dva standardy. Jelikož novější verze XEP dokumentu [13] se v době psaní této práce nachá- 238 zela ve stavu „experimental“, což znamená že ještě není schválena jako standard, je pouze 239 ve stavu návrhu. Proto bylo použito verze starší [16]. Jednoduše řečeno je *VCard* struktura, 240 která nese informace o uživateli jako je jméno, příjmení, e-mail, adresa bydliště i zaměst- 241 nání a další údaje. Data jsou dále zveřejňována na síti, z čehož vyplývá, že jsou dostupná 242 ostatním uživatelům. Vyplnění těchto osobních údajů je dobrovolné a tak se u některých 243 uživatelů nachází pouze přezdívka a JID, které jsou často předdefinovány automaticky. Ne- 244 dílnou součástí všech sociálních a komunikačních systému jsou malé fotografie, loga nebo 245 ikony, kterými se uživatelé prezentují. V síti Jabber tomu není jinak, a proto je samotný 246 obrázek zahrnut přímo do *VCard* v položce *photo*. Podrobnější informace o jeho nastavení 247 a přijímání je možné nalézt v *vCard-Based Avatars* [19], který jej definuje.

248 Díky základní podmínce XMPP protokolu (otevřenost) existuje mnoho různých aplikací 249 pomocí, kterých lze v síti Jabber komunikovat. S programy, používanými uživateli, úzce 250 souvisí další zde implementované rozšíření. Jedná se o realizaci *Software Version* dokumentu 251 [17], který se právě zabývá získáváním informací o samotných aplikacích. Je-li toto rozšíření 252 podporováno je díky němu možné zjistit jméno a verzi používané aplikace. Informace o 253 operačního systému často nejsou kvůli bezpečnosti ani vyplněny. Podrobnější informace o 254 softwarové výbavě klienta je možné zjistit pomocí XEP [7], o kterém již bylo dříve psáno v

255 odstavci zabývajícím se podporovanými vlastnostmi klientských aplikací.

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

265 Pro sdělení informací o stavu klienta není v základní verzi Jabberu mnoho. Pomocí
266 presence je možné „pouze“ prozradit zda je uživatel připraven komunikovat nebo je mo-
267 mentálně nedostupný a to v několika verzích lišících se délkou nepřítomnosti. Pokročilejší
268 nastavení statusu nabízí *User Mood* [20] a to ve formě sdělení současné nálady jako je
269 například radost. Další možné upřesnění činnosti uživatele jsou definovány v *User Activity*
270 [11], kde každá činnost je složena z povinné obecné kategorie a nepovinné, která informaci
271 upřesňuje. Příkladem může být *eating* a *having-a-snack* tj. uživatel jí, uživatel svačí.

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

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

283 Kapitola 3

284 Data mining

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

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

300 Terminologie

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

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

- 311 1. počet záznamů nebo-li objektů v databázi
- 312 2. počet polí nebo-li atributů objektů v databázi

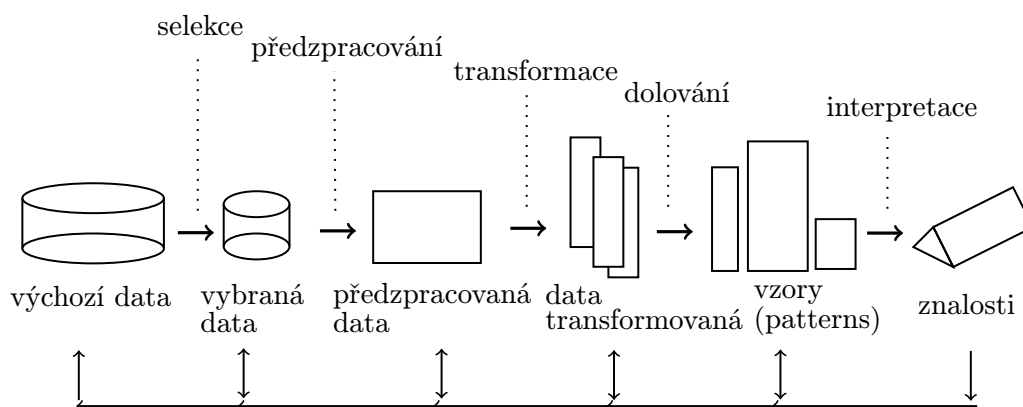
¹International Joint Conference on Artificial Intelligence

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

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

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

347 V této podsekcí jsou ve stručnosti vysvětleny základní nejdůležitější pojmy dále v práci
 348 využívány.

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

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

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

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

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

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

366 3.1 Metody dolování dat

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

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

- 375 • schopnosti reprezentace shluků (např. otázka lineární separability)
- 376 • srozumitelnosti nalezených znalostí pro uživatele (symbolické vs. subsymbolické me-
 377 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*.

408 Klasifikace

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

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

420 Modely

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

- 423 • Rozhodovací stromy
- 424 • Neuronové sítě
- 425 • Statistické metody
- 426 • Klasifikační pravidla
- 427 • Využití vzdálenosti
- 428 • a další ...

429 Predikce

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

437 Shlukování

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

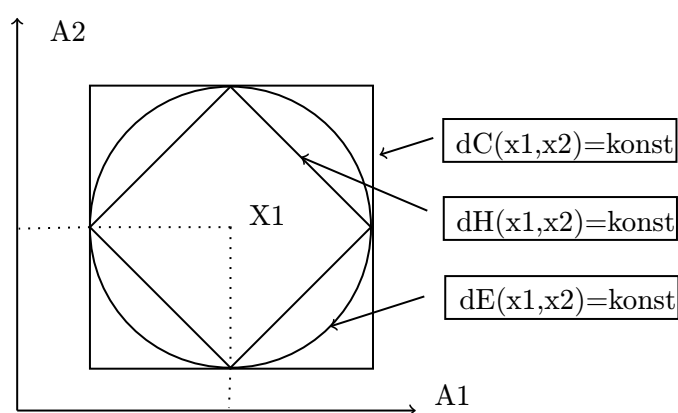
441 3.2 Shlukování

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

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

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

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



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

458 Metody založené na modelu

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

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

466 Metody hierarchické

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

⁴Self Organizing Oscillator Network

⁵Self-Organizing Map

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

474 Metody založené na mřížce

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

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

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

487 Metody založené na hustotě

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

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

501 Shlukování velkých dat

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

507 Zmíněné problémy shlukování velkých dat řeší dvě techniky *metoda transformace rysů*
508 *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

509 CLIQUE¹⁰.

510 Metody rozkladu

511 Metody rozkladu rozdělují datové prvky do několika podmnožin, nazývané shluky. Počet
512 shluků musí být znám před zahájením samotného procesu. Přiřazení do konkrétních tříd
513 je, podle [25], jednoznačné nebo probíhá na základě míry příslušnosti objektů do shluků.
514 Pro velký počet objektů, se kterými se pracuje, jsou využívány různé iterační optimalizace.

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

518 k -means

519 Shlukování pomocí algoritmu k -means je používáno pro data obsahující kvantitativní pro-
520 měnné a pro data, která nejsou příliš zašumělá. Základní proces je tvořen iterativním roz-
521 dělováním objektů do tříd na základě vzdáleností od jejich středů. Střed nebo-li centroid
522 shluku je vektor, jehož vzdálenost od součtu vzdáleností objektů v této třídě je minimální.
523 Pro výpočet vzdáleností mezi objekty samotnými nebo mezi objekty a středem je použita
524 euklidovská vzdálenost¹¹, která je vyobrazena na obrázku č. 3.2.

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

534 K -menas je algoritmus, kterým jsou přiřazovány objekty (vektory) x_n , kde $n = 1, \dots, N$,
535 do S_k , kde $k = 1, \dots, K$, shluků. V prvním kroku jsou určeny počáteční středy tříd, do
536 kterých se budou objekty shlukovat. Určení počátečních centroidů c_k probíhá například
537 náhodným výběrem K objektů nebo K prvních objektů souboru. Druhým krokem jsou
538 zkoumány jednotlivé vzdálenosti objektů x_n od počátečních středů c_j pomocí euklidovské
539 vzdálenosti. Na základě nejmenší zjištěné vzdálenosti mezi objektem a centroidem je ob-
540 jekt zařazen do shluku, kterému náleží právě tento střed. Třetím krokem tak jako u prvního
541 jsou hledány nové středy shluků, nyní ale již na matematickém základě. Je vypočítán na
542 základě průměrných jednotlivých hodnot objektů a uložen jako m -rozměrný vektor. Čtvr-
543 tým krokem se algoritmus dostává do konečné fáze, kdy mohou nastat dva možné případy.
544 Nově nalezené středy nejsou příliš vzdáleny od předchozích centroidů a proto je algoritmus
545 ukončen. Druhá častěji se vyskytující možnost iterativně provádí algoritmus od druhého
546 kroku dokud neplatí první možnost nebo dokud se objekty nepřestanou přemísťovat úplně.
547 Při popisu tohoto algoritmu bylo čerpáno z [25].

548 Díky jednoduchosti a relativní rychlosti je metoda k -means stále výrazně využívána.
549 Uplatnění nachází v široké škále oblastí jako je například biologie nebo počítačová grafika.

¹⁰CLustering In QUEst

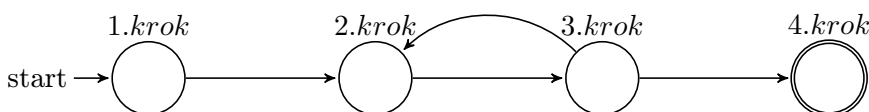
¹¹mean=střed, centroid je vektor průměrů

Algoritmus 3.1 Calculate $y = x^n$

Require: $n \geq 0 \vee x \neq 0$ **Ensure:** $y = x^n$

```
 $y \leftarrow 1$ 
if  $n < 0$  then
   $X \leftarrow 1/x$ 
   $N \leftarrow -n$ 
else
   $X \leftarrow x$ 
   $N \leftarrow n$ 
end if
while  $N \neq 0$  do
  if  $N$  is even then
     $X \leftarrow X \times X$ 
     $N \leftarrow N/2$ 
  else { $N$  is odd}
     $y \leftarrow y \times X$ 
     $N \leftarrow N - 1$ 
  end if
end while
```

550 Vzhledem k enormnímu počtu možného uspořádání nejsou výsledky vždy přesné, ale často
551 pouze přibližné.



Obrázek 3.3: My

552 3.3 Programy

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

563 RapidMiner

564 RapidMiner je, tak jak je popsán na oficiálních stránkách produktu [?], celosvětově nej-
565 používanější open-source systém pro dolování dat. Je možné jej používat jako samotnou
566 aplikaci nebo jej začlenit jako komponentu do vlastních výrobků. Pro zájemce je nabízen
567 také ve verzích pro firmy, které jsou rozdílné v poplatcích, podpoře pro zákazníka, záruce
568 a dalších balíčcích služeb zajišťující celkovou komplexnost a spolehlivost produktu.

569 Jak již většina podobných aplikací, je v současné době implementován v jazyce Java,
570 díky které nabízí flexibilní nejen grafické prostředí. K vybraným základním rysům toho
571 nástroje, tak jak jsou prezentovány firmou *Rapid-i*, patří: výkonné, přesto intuitivní grafické
572 uživatelské rozhraní pro návrh procesů, jednoduché řešení pro transformaci dat, kontrola
573 výsledků již při samotném návrhu a další. Nástrojem RapidMinerem je podporována široká
574 škála metod a algoritmů pro data minig, nejen vlastních a i z konkurenčního softwaru Weka.
575 V základní verzi určené pro veřejnost je k nalezení přes 100 procesů k modelování. Jsou
576 zde zastoupeny jak metody klasifikační a asociační, tak i metody shlukovací z nichž lze
577 jmenovat například DBSCAN, k -medoids a hlavně k -means.

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

585 Kapitola 4

586 Implementace

587 Obsahem čtvrté kapitoly je popis praktické části této práce. Jsou zde popsány jednotlivé
588 prvky, které byly použity jak pro získání dat, tak pro jejich následné uložení.

589 Cílem této práce je dolování dat z Jabberu. Jak již bylo dříve napsáno, Jabber je real-
590 time komunikační služba díky níž mohou její uživatelé komunikovat, informovat nebo sdílet
591 svůj status s jinými uživateli. Celá tato vzájemná komunikace skrývá nepřeborné množství
592 informací o klientech dané sítě. Všechna tato navzájem vyměňovaná nebo poskytnutá data
593 následně poslouží jako zdroj samotnému dolování. Pro jejich uskladnění je využita databáze,
594 jejichž strukturální návrh prezentuje obrázek č. 4.1.

595 Druhá část této kapitoly je zaměřena na robota, což je Jabber klient implementován
596 v jazyce C++ s konzolovým rozhraním. Robot v této práci hraje roli pasivního uživatele,
597 který informace pouze přijímá. Ve vybraných případech dokáže uživatele ve svém seznamu
598 kontaktů vyzvat k zaslání odpovědi s informacemi na požadovanou žádost. Struktura sa-
599 motného robota je popsána níže a reprezentována obrázkem č. ???. Na konci části o Jabber
600 robotovi bude zmíněno o knihovně *gloox* [?], kterou je zprostředkovány všechny náležitosti
601 Jabber komunikace.

602 4.1 Databáze

603 Data, která jsou sbírána robotem, jsou ukládána do objektově-relační databáze Postgre-
604 SQL¹. PostgreSQL nebo-li také Postgres byl použit ve verzi 8.4.7 a je provozován na ope-
605 račním systému Ubuntu Maverick verze 10.10.

606 Návrh databáze

607 Struktura databáze, do které je ukládána veškerá komunikace Jabber robota, využívá relační
608 model. Obrázkem č. 4.1 jsou prezentovány nejdůležitější části databáze. Celou strukturu
609 návrhu je možné nalézt v příloze ???. V jednotlivých částí návrhu databáze je počítáno s
610 jejím druhotným využitím obsahu, které bude popsáno v dalších částí této práce.

611 V době návrhu databáze zcela nebylo jasné, která data budou následně analyzována.
612 Z toho důvodu se struktura databáze snaží zachytit všechna „důležitá“ data. Předem ne-
613 bylo určen typ znalostí, kterých by se při dolování mělo dosáhnout. Za účelem neztratit
614 žádná data, je v návrhu databáze obsažena tabulka *xml*, která je nositelem obsah jak všech
615 přijatých, tak i odeslaných zpráv.

¹vyvinul se z projektu Ingres

616 **Transformace dat**

617 **4.2 Architektura**

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

620 **4.3 Robot**

621 –class diagram, pomoci nej popsati strukturu, mozne rozsireni

622 **gloox**

623 Gloox je stabilni Jabber/XMPP knihovna vydavana pod licenci GNU GPL. Je urcena
624 pro vyvoj klienta a komponent. Jelikož je psana v ANSCI C++ je mozne ji označit jako
625 multiplatformni ².

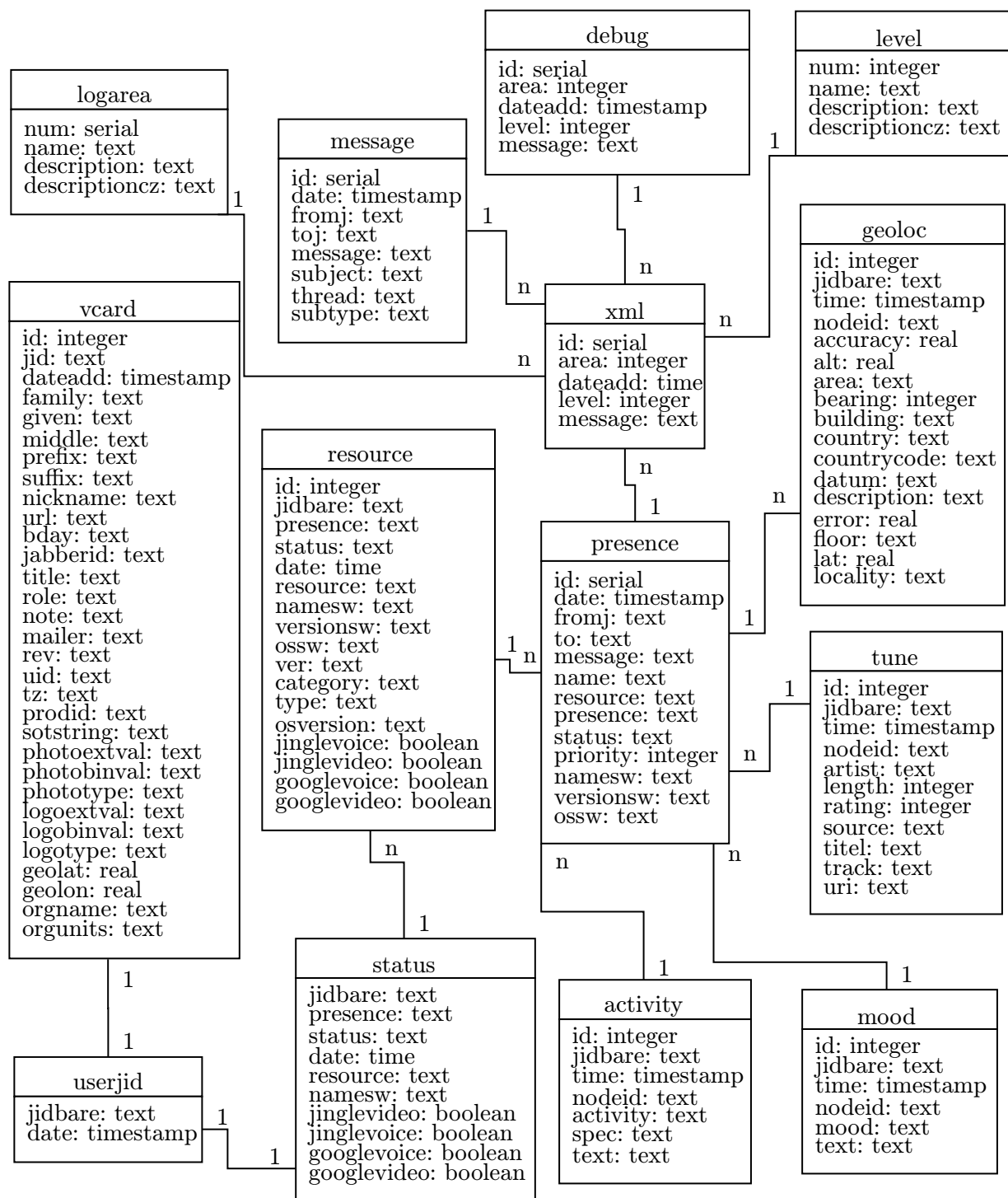
626 Pomoci knihovny gloox je psan bot v teto praci. Byla vybrana na zaklade požadavku
627 psani programu v jazyce C/C++ a operačním systémem Linux. V porovnání s jinými knihov-
628 nami pro jazyk C nebo C++ disponuje lepší podporou a dokumentací. Gloox plně podpo-
629 ruje standart XMPP Core [22] a z větší části i standard XMPP IM [23]. Dodatečně je plně
630 podporováno kolem 30 XEP standardů například vcard-temp[16] a další.

631 **Návrh bota**

632 **Transformace**

633 temporalni databaze

²Linux, Windows, Mac OS X, Symbian/Nokia S60, FreeBSD



Obrázek 4.1: Struktura databáze

⁶³⁴ Kapitola 5

⁶³⁵ Vyhodnocení výsledků

⁶³⁶ –charakter dat, jaka data jsem nasbiral –

⁶³⁷ Kapitola 6

⁶³⁸ Závěr

⁶³⁹

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.
26. 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. 26. 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.
26. 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.
26. 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. 26. dubna 2011].
URL <http://xmpp.org/extensions/xep-0060.html>

- 670 [13] Mizzi, S.; Saint-Andre, P.: XEP-0292: vCard4 Over XMPP. [online], 02-26-2008, [cit.
671 26. dubna 2011].
672 URL <http://xmpp.org/extensions/xep-0292.html>
- 673 [14] Moore, D.; Wright, W.: *Jabber developer's handbook*. Indianapolis: Sams Publishing,
674 první vydání, 2004, 487 s., iISBN 06-723-2536-5.
- 675 [15] Nemrava, M.; Pospíšil, J.: Dolování dat a jeho aplikace. [online], 2006, [cit. 26. dubna
676 2011].
677 URL http://www.spatial.cs.umn.edu/paper_ps/dmchap.pdf
- 678 [16] Saint-Andre, P.: XEP-0054: vcard-temp. [online], 07-16-2008, [cit. 26. dubna 2011].
679 URL <http://xmpp.org/extensions/xep-0054.html>
- 680 [17] Saint-Andre, P.: XEP-0092: Software Version. [online], 02-15-2007, [cit. 26. dubna
681 2011].
682 URL <http://xmpp.org/extensions/xep-0092.html>
- 683 [18] Saint-Andre, P.: XEP-0118: User Tune. [online], 30-01-2008, [cit. 26. dubna 2011].
684 URL <http://xmpp.org/extensions/xep-0118.html>
- 685 [19] Saint-Andre, P.: XEP-0153: vCard-Based Avatars. [online], 16-08-2006, [cit.
686 26. dubna 2011].
687 URL <http://xmpp.org/extensions/xep-0153.html>
- 688 [20] Saint-Andre, P.; Meijer, R.: XEP-0107: User Mood. [online], 29-10-2008, [cit.
689 26. dubna 2011].
690 URL <http://xmpp.org/extensions/xep-0107.html>
- 691 [21] Saint-Andre, P.; Smith, K.: XEP-0163: Personal Eventing Protocol. [online],
692 12-07-2010, [cit. 26. dubna 2011].
693 URL <http://xmpp.org/extensions/xep-0163.html>
- 694 [22] Saint-André, P.: Extensible Messaging and Presence Protocol (XMPP): Core.
695 [online], 10-2004, [cit. 26. dubna 2011].
696 URL <http://tools.ietf.org/html/rfc3920>
- 697 [23] Saint-André, P.: Extensible Messaging and Presence Protocol (XMPP): Instant
698 Messaging and Presence. [online], 10-2004, [cit. 26. dubna 2011].
699 URL <http://tools.ietf.org/html/rfc3921>
- 700 [24] Saint-André, P.; Smith, K.; Troncon, R.: *XMPP : the definitive guide : building*
701 *real-time applications with jabber technologies*. Sebastopol: O'Reilly, první vydání,
702 2009, 287 s., iISBN 978-059-6521-264.
- 703 [25] Řezánková, H.; Húsek, D.; Snášel, V.: *Shluková analýza dat*. Praha: Professional
704 Publishing, druhé vydání, 2009, 218 s., iISBN 978-808-6946-818.

⁷⁰⁵ **Příloha A**

⁷⁰⁶ **Obsah CD**

⁷⁰⁷ **Příloha B**

⁷⁰⁸ **Manual**

⁷⁰⁹ **Příloha C**

⁷¹⁰ **Konfigurační soubor**

711 Příloha D

712 Slovník zkratek

713 **DNS** — Domain Name System

714 **GPG** — dkshckdsjvlsdjvodsvjdfokj

715 **IM služby** — dkshckdsjvlsdjvodsvjdfokj

716 **IP** — Internet Protocol

717 **JEP** — dkshckdsjvlsdjvodsvjdfokj

718 **JID** — dkshckdsjvlsdjvodsvjdfokj

719 **SASL** — dkshckdsjvlsdjvodsvjdfokj

720 **TCP** — dkshckdsjvlsdjvodsvjdfokj

721 **TLS** — dkshckdsjvlsdjvodsvjdfokj

722 **WWW** — dkshckdsjvlsdjvodsvjdfokj

723 **XEP** — dkshckdsjvlsdjvodsvjdfokj

724 **XML** — dkshckdsjvlsdjvodsvjdfokj

725 **XMPP** — dkshckdsjvlsdjvodsvjdfokj

726 **e-mail** — dkshckdsjvlsdjvodsvjdfokj

727 **jabber** — dkshckdsjvlsdjvodsvjdfokj

728 **klient** — dkshckdsjvlsdjvodsvjdfokj

729 **presence** — dkshckdsjvlsdjvodsvjdfokj

730 **server** — dkshckdsjvlsdjvodsvjdfokj

731 **stanza** — dkshckdsjvlsdjvodsvjdfokj

732 **vCard** — dkshckdsjvlsdjvodsvjdfokj

733 **ID3v1** — dkshckdsjvlsdjvodsvjdfokj

734 **ID3v2** — dkshckdsjvlsdjvodsvjdfokj

735 Příloha E

736 Stanza - základní schéma

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

741 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*.

742 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*.

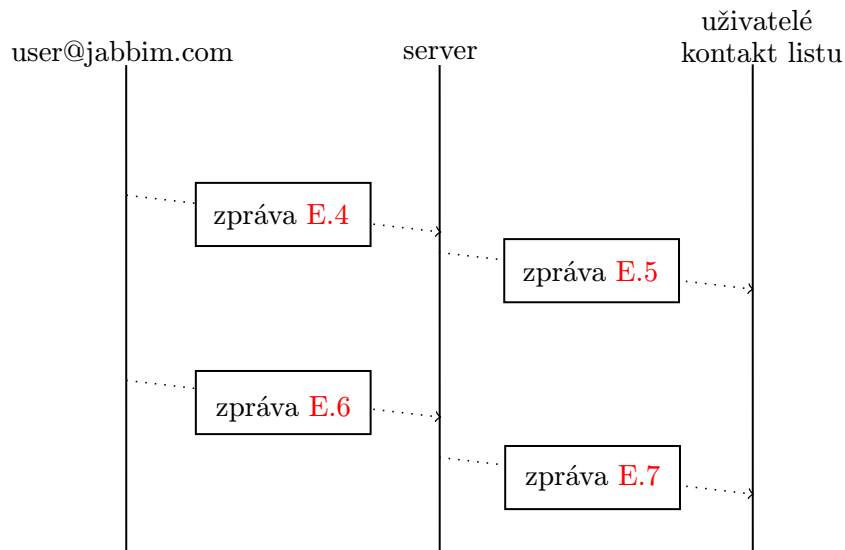
E.3 Presence

```
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*.

744 E.4 Přehled průběhu rozšíření

745 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*.

746

```

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.

747 Server obdrží informace od klienta *user* zprávu o přehrávací hudbě. Pomocí elementu
748 *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*.

749 Zpráva o přehrávané hudbě je také přeposlána všem otevřeným spojením uživatele *user*,
750 ukázáno na příkladě E.6.
751

```
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*.

752 Přestane-li uživatel *user* poslouchat/vysílat informace o přehrávané hudbě, provede to
753 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.

754 Server informuje všechny účastníky odběru zprávou, která má položku *tune* prázdnou.
755 Tak jak to prezentuje příklad E.8.
756

```
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*.

757 Příloha F

758 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 | | | | | | | | | | |
| Ayttm | | | | | | | | | | |
| 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.