

**GREGORY MORO PUPPI WANDERLEY**

**FolksDialogue: Um Método para o Aprendizado  
Automático de Folksonomias a partir de Diálogo  
Orientado à Tarefa em Português do Brasil**

Dissertação de Mestrado apresentada ao  
Programa de Pós-Graduação em Informática da  
Pontifícia Universidade Católica do Paraná  
como requisito parcial para obtenção do título de  
Mestre em Informática.

**CURITIBA**

**2015**

**GREGORY MORO PUPPI WANDERLEY**

**FolksDialogue: Um método para o Aprendizado Automático de Folksonomias a partir de Diálogo Orientado à Tarefa em Português do Brasil**

Dissertação de Mestrado apresentada ao Programa de Pós-Graduação em Informática da Pontifícia Universidade Católica do Paraná como requisito parcial para obtenção do título de Mestre em Informática.

Área de Concentração: *Ciência da Computação*

Orientador: Prof. Dr. Emerson Cabrera Paraiso

**CURITIBA**

**2015**

Wanderley, Gregory Moro Puppi

W245f FolksDialogue : um método para o aprendizado automático de folksonomias  
2015 a partir de diálogo orientado à tarefa em português do Brasil / Gregory Moro  
Puppi Wanderley ; orientador, Emerson Cabrera Paraiso. – 2015.  
xiii, 121 f. : il. ; 30 cm

Dissertação (mestrado) – Pontifícia Universidade Católica do Paraná,  
Curitiba, 2015

Bibliografia: f. [112-121]

1. Computação. 2. Programas de aprendizado. 3. Diálogos. 4. Indexação.  
I. Paraiso, Emerson Cabrera. II. Pontifícia Universidade Católica do Paraná.  
Programa de Pós-Graduação em Informática. III. Título.

CDD 20. ed. – 004


Dados da Catalogação na Publicação  
Pontifícia Universidade Católica do Paraná  
Sistema Integrado de Bibliotecas – SIBI/PUCPR  
Biblioteca Central

**ATA DE DEFESA DE DISSERTAÇÃO DE MESTRADO  
PROGRAMA DE PÓS-GRADUAÇÃO EM INFORMÁTICA**

**DEFESA DE DISSERTAÇÃO DE MESTRADO Nº 02/2015**

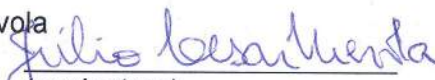
Aos 05 dias do mês de Março de 2015 realizou-se a sessão pública de Defesa da Dissertação “ **FolksDialogue: Um Método para o Aprendizado Automático de Folksonomias a partir de Diálogo Orientado à Tarefa em Português do Brasil**” apresentado pelo aluno **Gregory Moro Puppi Wanderley**, como requisito parcial para a obtenção do título de Mestre em Informática, perante uma Banca Examinadora composta pelos seguintes membros:

Prof. Dr. Emerson Cabrera Paraiso  
PUCPR (Orientador)

  
(assinatura)


APROVADO  
(Aprov/Reprov)

Prof. Dr. Júlio César Nievola  
PUCPR

  
(assinatura)


APROVADO  
(Aprov/Reprov)

Prof. Dr. Cesar Augusto Tacla  
UTFPR

  
(assinatura)

Aprovado  
(Aprov/Reprov)

Prof. Dr. Sylvio Barbon Júnior  
UEL

  
(assinatura)

APROVADO  
(Aprov/Reprov)

Conforme as normas regimentais do PPGIa e da PUCPR, o trabalho apresentado foi considerado APROVADO (aprovado/reprovado), segundo avaliação da maioria dos membros desta Banca Examinadora. Este resultado está condicionado ao cumprimento integral das solicitações da Banca Examinadora registradas no Livro de Defesas do programa.

  
Prof.ª Dr.ª Andreia Malucelli.

Coordenadora do Programa de Pós-Graduação em Informática.



*Dedico esta dissertação aos meus pais Jérvys e Luzimeri,  
e à minha irmã Madeleine.*

*“Somente se aproxima da perfeição  
quem a procura com constância,  
sabedoria e, sobretudo, humildade.”  
(Jigoro Kano)*

## Agradecimentos

Para que o desenvolvimento deste trabalho fosse realizado e seus objetivos concluídos com sucesso, foram necessários o apoio, o auxílio e a contribuição de diversas pessoas. Aproveito para deixar meus agradecimentos a todos que participaram ou contribuíram de algum modo, diretamente ou indiretamente, com esta pesquisa.

Em primeiro lugar gostaria de agradecer à minha família, especialmente aos meus pais, que sempre me apoiaram incondicionalmente em toda a trajetória de minha vida. Agradeço pela base, pelos valores, por todos os ensinamentos e oportunidades que me concederam. À minha irmã, por estar sempre ao meu lado e apta a me auxiliar. A vocês, minha eterna gratidão.

Ao meu orientador, professor Emerson, pela oportunidade, por ter acreditado em mim, pelo conhecimento fornecido, por ter me direcionado nos momentos mais críticos, e pela paciência durante todo o período de orientação. A você, toda a minha admiração e reconhecimento.

Aos professores, Julio Cesar Nievola e Cesar Augusto Tacla, pela participação em minha banca de qualificação e por todos os conselhos transmitidos. Agradeço vocês por todo o auxílio e disposição.

A todos os colegas e amigos do PPGIa e do laboratório LASIN, em especial, Ednilson, Mariza, Irapuru, Andréia, Franciele, Elias e Patrícia. À Cheila, secretária do PPGIa. Obrigado pela amizade, ajuda e presteza.

À CAPES por todo o suporte e financiamento prestado durante o projeto de mestrado. Agradeço por tornar real este sonho.

Por fim, a todos aqueles que posso ter esquecido de mencionar neste momento. Meu mais profundo agradecimento.

# Sumário

<b>LISTA DE FIGURAS .....</b>	<b>VIII</b>
<b>LISTA DE QUADROS .....</b>	<b>IX</b>
<b>LISTA DE TABELAS.....</b>	<b>X</b>
<b>LISTA DE ABREVIATURAS.....</b>	<b>XI</b>
<b>RESUMO.....</b>	<b>XII</b>
<b>ABSTRACT .....</b>	<b>XIII</b>
<b>CAPÍTULO 1 INTRODUÇÃO .....</b>	<b>14</b>
1.1. Motivação.....	17
1.2. Objetivos.....	19
1.3. Hipóteses de Trabalho .....	19
1.4. Contribuição Científica.....	19
1.5. Escopo.....	20
1.6. Organização do Documento.....	20
<b>CAPÍTULO 2 FOLKSONOMIAS .....</b>	<b>22</b>
2.1. Folksonomias .....	22
2.2. Grafos .....	26
2.2.1. Representação de Folksonomias por Grafos .....	28
2.2.2. Mundo-Pequeno.....	30
2.3. Conclusão .....	33
<b>CAPÍTULO 3 DIÁLOGOS .....</b>	<b>34</b>
3.1. Diálogos .....	34
3.2. Interpretador de Linguagem Natural .....	38
3.3. Conclusão .....	40
<b>CAPÍTULO 4 TRABALHOS RELACIONADOS .....</b>	<b>41</b>
4.1. Aprendizado de Folksonomias.....	42
4.2. Avaliação de Folksonomias .....	44
4.3. Detecção de Tendências em Folksonomias .....	45
4.4. Visão Geral dos Trabalhos .....	47
4.5. Análise dos Trabalhos Relacionados e Conclusão.....	54
<b>CAPÍTULO 5 UM MÉTODO PARA O APRENDIZADO AUTOMÁTICO DE FOLKSONOMIAS A PARTIR DE DIÁLOGOS.....</b>	<b>56</b>
5.1. Definição Formal de Folksonomias obtidas a partir de Diálogos .....	57
5.2. O Método FolksDialogue .....	60
5.2.1. Etapas da Atividade de Pré-processar .....	65
5.2.2. Etapas da Atividade de Aprender .....	68
5.3. Conclusão .....	77
<b>CAPÍTULO 6 PROCEDIMENTOS METODOLÓGICOS.....</b>	<b>78</b>
6.1. Corpus de Diálogos.....	78
6.2. Implementação do Protótipo Computacional .....	79
6.2.1. Implementação da Atividade de Pré-processar .....	80
6.2.2. Implementação da Atividade de Aprender .....	81
6.3. Avaliação.....	88
6.3.1. Avaliação de Característica .....	88
6.3.2. Avaliação de Teste de Domínio.....	91
6.4. Abordagem de Detecção de Tendências .....	93
6.6. Conclusão .....	98
<b>CAPÍTULO 7 RESULTADOS OBTIDOS .....</b>	<b>99</b>



<b>7.1. Corpus de Entrada .....</b>	<b>99</b>
<b>7.2. Resultados da Avaliação de Característica .....</b>	<b>101</b>
<b>7.3. Resultados da Avaliação de Teste de Domínio .....</b>	<b>102</b>
<b>7.4. Resultados da Abordagem de Detecção de Tendências .....</b>	<b>105</b>
<b>7.5. Conclusão .....</b>	<b>107</b>
<b>CAPÍTULO 8 CONCLUSÕES E TRABALHOS FUTUROS.....</b>	<b>108</b>
<b>REFERÊNCIAS BIBLIOGRÁFICAS .....</b>	<b>112</b>

## Lista de Figuras

<b>Figura 1</b> - Um exemplo de folksonomia (Fonte: adaptado de DATTOLO e PITASSI (2012)). .....	24
<b>Figura 2</b> - Exemplo de relacionamento entre rótulos de uma folksonomia (adaptado de PLANGPRASOPCHOK e LERMAN (2009)). .....	25
<b>Figura 3</b> - Exemplo da representação geométrica dos grafos. .....	26
<b>Figura 4</b> - Exemplo de grafo rotulado. .....	27
<b>Figura 5</b> - Um grafo tripartido, com três conjuntos de vértices: A, B e C (Fonte: adaptado de CHOJNACKI e KLOPOTEK (2010)). .....	27
<b>Figura 6</b> - Um exemplo de folksonomia. .....	28
<b>Figura 7</b> - Um grafo tripartido que representa um trecho de uma folksonomia. .....	29
<b>Figura 8</b> - Um trecho de uma folksonomia com os relacionamentos entre tags. .....	30
<b>Figura 9</b> - Exemplo de diálogo orientado a tarefa. .....	38
<b>Figura 10</b> - Frame para mostra horários de linha de ônibus. .....	39
<b>Figura 11</b> - Gramática para extração de significado. .....	39
<b>Figura 12</b> - Modelo quadripartido que representa a folksonomia obtida a partir de diálogos. .....	59
<b>Figura 13</b> - Diagrama de atividades representando o fluxo do método. .....	61
<b>Figura 14</b> - Etapas da atividade de Pré-processar que compõe o método proposto. .....	62
<b>Figura 15</b> - Etapas da atividade de Aprender que compõe o método proposto. .....	64
<b>Figura 16</b> - Folksonomia aprendida pelo método FolksDialogue. .....	76
<b>Figura 17</b> - Personomia do atendente $a_2$ , extraída da folksonomia aprendida pelo método FolksDialogue. .....	77
<b>Figura 18</b> - Diagrama simbolizando o fluxo da abordagem proposta para detecção de tendências. .....	95
<b>Figura 19</b> - Detalhes do processo que realiza a extração dos Assuntos Abordados. .....	96

## Lista de Quadros

<b>Quadro 1</b> - Bases Consultadas e Strings de Busca (Fonte: Autor). .....	41
<b>Quadro 2</b> - Comparação entre os Trabalhos Relacionados (Fonte: Adaptado de (FERREIRA, 2013)). .....	48
<b>Quadro 3</b> - Exemplos de Diálogos Orientados à Tarefas (Fonte: Autor). .....	62
<b>Quadro 4</b> - Resultado da Extração dos Enunciados dos Atendentes (Fonte: Autor). .....	66
<b>Quadro 5</b> - Resultado da Extração dos Substantivos dos Enunciados dos Atendentes (Fonte: Autor). .....	67
<b>Quadro 6</b> - Substantivos extraídos dos enunciados dos atendentes (Fonte: Autor). .....	67
<b>Quadro 7</b> - Conjunto de substantivos resultante da remoção de duplicados (Fonte: Autor). .....	67
<b>Quadro 8</b> - Ranqueamento dos Substantivos (Fonte: Autor). .....	69
<b>Quadro 9</b> - Enunciados dos Atendentes que são os Recursos da folksonomia Proposta (Fonte: Autor). .....	71
<b>Quadro 10</b> - Pares de Rótulos: com suas Frequências e seus Pesos $w$ (Fonte: Autor). .....	73
<b>Quadro 11</b> - Exemplos de Diálogos do Corpus Utilizado nesta Pesquisa. ....	79
<b>Quadro 12</b> - Algoritmo para criação das personomias dos Atendentes (Fonte: Autor). .....	85
<b>Quadro 13</b> - Algoritmo para Ligação dos Rótulos das personomias dos Atendentes (Fonte: Autor). .....	88
<b>Quadro 14</b> – Exemplo do processo de interpretação de enunciados de diálogos pela folksonomia aprendida (Fonte: Autor). .....	92
<b>Quadro 15</b> – Exemplos de Diálogos Utilizados como Corpus de Entrada no Método Proposto. ....	100
<b>Quadro 16</b> – Exemplo de Assuntos Abordados que se tornaram tendências. ....	106

## Lista de Tabelas

<b>Tabela 1</b> - Parâmetros da Folksonomia Apreendida para a Avaliação de Característica. ....	101
<b>Tabela 2</b> – Resultados da avaliação de característica. ....	102
<b>Tabela 3</b> – Assuntos Abordados extraídos por cada folksonomia. ....	105

## Lista de Abreviaturas

PLN	<i>Processamento de Linguagem Natural</i>
URL	<i>Uniform Resource Locator</i>
OWL	<i>Web Ontology Language</i>

## Resumo

Sistemas de diálogo têm como intuito facilitar a interação entre seres humanos e computadores, tornando-a mais natural. Esses sistemas permitem que um humano interaja com um computador utilizando o processo de interpretação de linguagem natural. Um dos componentes fundamentais desse processo em sistemas de diálogos é o Modelo Conceitual. O Modelo Conceitual representa um dado domínio e sua especificação se dá por diversas formas de representação do conhecimento. Nesta pesquisa propõe-se representar o Modelo Conceitual através de folksonomias, as quais são estruturas de representação do conhecimento que emergem do processo de rotulação, em sistemas de rotulação colaborativa. O processo de rotulação corresponde às atribuições de rótulos a recursos, por usuários. Desse modo, folksonomias são compostas por usuários, rótulos e recursos. Recursos podem ser quaisquer objetos que os usuários tenham interesse em rotular, como, fotos, vídeos e músicas. A coleção das rotulações geradas por um dado usuário corresponde a sua personomia. Assim, pode-se dizer que uma folksonomia é a reunião das personomias de todos os usuários que participaram da criação dos rótulos que descrevem um domínio qualquer. Uma das características interessantes das folksonomias é a sua dimensão social (usuários), a qual também é um aspecto presente em diálogos, resultante da interação entre os humanos. Nesta dissertação é descrito um método que realiza o aprendizado automático de folksonomias a partir de diálogo orientado à tarefa em português do Brasil. As folksonomias obtidas pelo método podem ser úteis no processo de interpretação de enunciados de diálogos, indicando se eles pertencem ou não aos domínios que elas representam. Além disso, pode ser possível realizar a detecção de tendências nelas, como por exemplo, verificando-se quais são os assuntos abordados em um dado domínio de diálogos, num determinado intervalo de tempo. Experimentos envolvendo uma base real de diálogos orientados a tarefas mostraram que as folksonomias, aprendidas pelo FolksDialogue, podem interpretar enunciados com uma acurácia de 72,32%. Um experimento confirmou que também é possível determinar tópicos abordados pelos interlocutores nos diálogos.

**Palavras-Chave:** Modelo Conceitual, Folksonomias, Diálogos, Aprendizado Automático.

## Abstract

Dialogue systems are intended to facilitate the interaction between humans being and computers, making it more natural. A dialogue system allows a human to interact with a computer, through the natural language processing. One of the main components of the natural language processing in dialogue systems is the Conceptual Model. The Conceptual Model represents a domain and its specification is given through various forms of knowledge representation. In this research, we propose to represent the Conceptual Model through folksonomies, which are structures of knowledge representation that emerges from the tagging process, in collaborative tagging systems. The tagging process corresponds to the assignments of tags to resources, by users. Thus, folksonomies are composed by users, tags and resources. Resources can be any objects, which users have interest in tagging, such as photos, videos and music. The collection of all assignments generated by a user is called his personomy. Based on this, we can say that a folksonomy is the union of all personomies of all users that have participated in the tagging process of some domain. One of the interesting features of folksonomies, is its social dimension (users), which is also presented in dialogues, resulting from the interaction process among humans being. This research describes a method that performs the automatic learning of folksonomies from task-oriented dialogues in Brazilian Portuguese. The folksonomies generated by the method, can be useful in the interpretation of dialogue utterances, indicating whether the utterances belong or not to the domains that the folksonomies represent. Besides that, it can also be possible to detect some trends with the obtained folksonomies, such as verifying what are the discussed topics in a given domain of dialogues in a given time interval. Experiments involving a real-world task-oriented dialogue corpus showed that using our method, learned folksonomies can interpret utterances with an accuracy of 72.32%. Moreover, an experiment confirmed that it is possible to determine topics addressed by interlocutors in dialogues.

**Keywords:** Conceptual Model, Folksonomies, Dialogues, Automatic Learning.

# Capítulo 1

## Introdução

No intuito de facilitar a interação entre humanos e computadores, diminuindo a necessidade de conhecimento prévio da interface computacional, e tornando esta interação mais natural, sistemas de diálogo ou interfaces conversacionais podem ser utilizados (ERIKSSON, 1999). Um sistema de diálogo permite que um humano interaja com o computador, como se a máquina fosse outro ser humano. Sistemas de diálogos podem ser encontrados numa larga gama de aplicações, como por exemplo, em *websites* (respondendo as dúvidas de clientes sobre produtos e serviços), em suporte técnico (respondendo ou diagnosticando problemas), em aprendizagem (dando conselhos enquanto o usuário aprende sobre um dado assunto) (LESTER et al., 2004).

Uma das áreas de pesquisa em que sistemas de diálogo se apoiam para interpretar os enunciados apresentados pelos usuários é a do Processamento de Linguagem Natural (PLN). O processo de interpretação de linguagem natural passa por uma série de etapas que envolvem a análise do conhecimento morfológico, sintático, semântico e pragmático (JURAFSKY e MARTIN, 2008).

Existe também como parte integrante do processo de interpretação de linguagem natural, a interação de todas essas etapas com repositórios que armazenam informações sobre o domínio de que se trata a aplicação. Um exemplo dessas informações é o Modelo Conceitual.

De acordo com (GUIZZARDI, 2005), um *modelo* é uma abstração de uma parte da realidade que está presente somente na mente das pessoas. O modelo é formado por um conjunto de *conceitos* usados para articular abstrações do estado das coisas num dado domínio (GUIZZARDI, 2005). Esse modelo de conceitos ou Modelo Conceitual fornece conceitos



específicos, ou seja, conceitos referentes a um domínio particular do conhecimento ou relacionados a tarefas específicas (DI FELIPPO e DIAS DA SILVA, 2006). A concretização do Modelo Conceitual é denominada de *especificação do Modelo Conceitual*. O intuito da especificação é tornar o Modelo Conceitual apto a ser documentado e analisado. Além disso, a especificação do Modelo Conceitual é usada para dar suporte ao entendimento entre as partes interessadas num dado domínio (GUIZZARDI, 2005). Assim, a especificação do Modelo Conceitual, intitulada de *Modelo Conceitual*, é um artefato concreto que permite aos atores envolvidos no processo de construção do modelo compreender o domínio (FERREIRA, 2013).

O processo de transformação de modelos abstratos e de seus conceitos para modelos concretos envolve a *Aquisição do Conhecimento* (GOMÉZ-PÉREZ et al., 2004), que fornece aos atores envolvidos conhecimentos necessários para a especificação do Modelo Conceitual (FERREIRA, 2013).

O Modelo Conceitual pode ser representado por diversas formas, como ontologias (GUARINO et al., 2009) e frames (MINSKY, 1975). Embora existam trabalhos que realizem a obtenção automática dessas estruturas (HAZMAN et al., 2011), normalmente a Aquisição do Conhecimento é realizada por engenheiros de conhecimento (LIU e GRUEN, 2008) (WEAVER et al., 1988), em oposição a usuários comuns. Em geral para a obtenção do conhecimento, os engenheiros de conhecimento utilizam técnicas como entrevistas, *brainstormings* (FARZANEH et al., 2013) ou ferramentas de apoio a aquisição (BOZ JR et al., 2011). Essas técnicas foram criadas para Sistemas a Base de Conhecimento (como os sistemas especialistas), os quais são dedicados a aplicações particulares com fontes de informações e quantidade de atores restrita. No entanto, hoje em dia as fontes de informação são maiores (WANG et al., 2006), como por exemplo a *Web*, e o número de atores envolvidos (engenheiros de conhecimento, especialistas e usuários) também pode ser maior (TEMPICH et al., 2005).

Quando a especificação do Modelo Conceitual do domínio envolve um grande número de atores no processo, as técnicas de aquisição de conhecimento que exigem a interação direta entre engenheiro de conhecimento e especialistas como entrevistas e *brainstormings* podem demandar muito tempo, custo e esforço, devido à dificuldade em se atingir consenso entre os atores. Atingir consenso para especificar os conceitos de um domínio não é uma tarefa simples, principalmente quando o número de pessoas envolvidas cresce, pois aumentam as

divergências, assim como o número de interações para resolvê-las. Segundo a semiótica, um dos motivos desta complexidade é devido a ser muito comum o fato de duas pessoas, ao se comunicarem gerarem interpretações diferentes para uma mesma entidade ou representação (ECO, 1976) (FERREIRA, 2013).

Segundo (FERREIRA, 2013), na *Web* um dos serviços que ganhou popularidade entre os usuários nos últimos anos, foram os sistemas baseados em *tagging* ou rotulação colaborativa. Em termos de especificação de modelos conceituais, o que se destaca nesses sistemas é uma espécie de “conhecimento coletivo” que emerge da contribuição individual de cada usuário (GRUBER, 2007). Em outras palavras, o conhecimento é gerado pelos próprios usuários, sendo independente de um grupo de especialistas que necessitam encontrar a igualdade de opiniões e/ou pensamento de como representar um domínio, algo comum nas ontologias (GRUBER, 1993). A estrutura de representação do conhecimento que emerge do processo de *tagging* ou rotulação, ou seja, da interação dos usuários com o sistema, é denominada de *folksonomia*, junção de “*folk*” e “*taxonomia*” (VAN DER WAL, 2004).

Sistemas baseados em rotulação colaborativa são aplicações sociais que permitem aos seus usuários atribuírem *tags* ou rótulos a recursos, como *URLs*, fotos, vídeos, dentre outros. Deste modo, constata-se que tanto o processo de rotulação, quanto as folksonomias que emergem dele são compostos de usuários, rótulos e recursos. Em um sistema de rotulação, a *personomia* de um usuário é o conjunto de todos os rótulos e os recursos que foram rotulados por ele (SILVA et al., 2012). Assim, pode-se inferir que uma folksonomia é o conjunto das personomias de todos os usuários que participaram do processo de rotulação de um domínio qualquer.

De acordo com (GUPTA et al., 2010), algumas aplicações em que o processo de rotulação colaborativa tem sido útil são: indexação, busca, melhorar a navegação, dentre outros. Em termos de indexação, rótulos podem acelerar esse processo em *sites* ou páginas da internet. Para busca, o trabalho de (BAO et al., 2007) observou que o processo de rotulação pode trazer benefícios em dois aspectos: rótulos relacionados a um *site* são normalmente bons resumos do que eles representam, e o número de rotulações num dado site indica a popularidade deles, servindo com isso de auxílio para algoritmos de ranqueamento de busca, como o *PageRank* (BRIN e PAGE, 1998) do Google<sup>1</sup>. No âmbito de melhorar a navegação, o trabalho de (ZUBIAGA, 2009) sugeriu maneiras alternativas de navegação através do uso de

---

<sup>1</sup> Maiores informações em: <<http://www.google.com/>>

rótulos. A primeira é denominada “*Pivot-browsing*”, e significa mover-se num espaço de informação através da escolha de um ponto de referência, por exemplo, indo-se até um rótulo é possível verificar quais rótulos estão relacionados a ele. A segunda chamada de “*Popularity-driven navigation*”, significa que algumas vezes os usuários podem querer recuperar somente os documentos que são rotulados por um determinado rótulo popular entre eles. E a terceira, “*Filtering*”, tem como intuito recuperar apenas os documentos que possuem um dado rótulo, excluindo os demais.

### 1.1. Motivação

O Modelo Conceitual descreve um dado domínio e é parte fundamental do processo de interpretação dos enunciados de um sistema de diálogo. A especificação ou construção do Modelo Conceitual pode ser feita por diversas formas de representação do conhecimento, como ontologias e folksonomias. Geralmente a especificação através de ontologias requer um grupo de especialistas (WELLER, 2007) (PAULSEN et al., 2007) que necessitam encontrar a igualdade ou um consenso de opiniões de como representar um domínio (GRUBER, 1993). No entanto, quando as fontes de informação são maiores, como na *Web*, e também quando o número de atores (especialistas e usuários) envolvidos no processo é maior, a obtenção de ontologias pode se tornar custosa e demorada devido à dificuldade em se atingir o consenso.

Ontologias são formas de representação do conhecimento complexas (WELLER, 2007), nas quais conceitos, instâncias, atributos e relacionamentos são modelados. Assim, devido ao fato das ontologias modelarem um domínio de forma rigorosa ou formal, elas são difíceis de serem obtidas. Outro fator envolvido com as ontologias é a manutenção ou atualização de suas estruturas com a inserção de novas informações, principalmente em ambientes dinâmicos (ECHARTE et al., 2004), algo que tende a ser custoso. Segundo (GOMÉZ-PÉREZ et al., 2004), as ontologias necessitam ainda de considerações com as linguagens a serem utilizadas para as suas construções, como por exemplo, os três tipos de camadas da linguagem *OWL*<sup>2</sup>.

Por outro lado, folksonomias são estruturas de representação do conhecimento que emergem do processo de *tagging* ou rotulação em sistemas de rotulação colaborativa (XIAO

---

<sup>2</sup> Maiores informações em: <<http://www.w3.org/TR/owl-features/>>

et al., 2010) (CATTUTO et al., 2007) (PETERS, 2009). O processo de rotulação corresponde às atribuições de rótulos a recursos por usuários. Deste modo, as folksonomias são compostas por usuários, *tags* ou rótulos e recursos. Recursos podem ser quaisquer objetos que os usuários tenham interesse em rotular, como por exemplo, fotos, vídeos e músicas. Em comparação com as ontologias que são estruturas complexas, folksonomias são formas de representação mais simples de se implementar e utilizar (ECHARTE et al., 2004). Segundo (HOTH0 et al., 2006a), um dos benefícios do processo de rotulação é que os usuários não precisam de experiência ou habilidades específicas para participar, ou seja, as folksonomias que emergem não necessitam ser construídas por engenheiros de conhecimento. Além disso, diferente de ontologias que possuem um vocabulário controlado obtido a partir do consenso entre especialistas, folksonomias refletem diretamente o vocabulário de usuários comuns, pois são eles mesmos quem rotulam os recursos (QUINTARELLI, 2005). Assim, as folksonomias em oposição às ontologias não sofrem com o grande volume de informação, nem com a necessidade de consenso. No entanto, vale ressaltar que as folksonomias possuem algumas limitações, como por exemplo, o fato de serem compostas apenas por usuários, rótulos e recursos. Sendo que tais rótulos são palavras-chave, as quais não possuem atributos e/ou relação hierárquica entre si, e nem definição de seus significados.

Uma das características interessantes das folksonomias é a sua dimensão social (usuários), a qual também é um aspecto presente em diálogos, resultante do processo de interação entre os seres humanos. Deste modo, esta pesquisa propõe fazer o aprendizado de folksonomias a partir de diálogos e posteriormente utilizar as estruturas obtidas para representar o Modelo Conceitual, útil na interpretação de enunciados de um dado domínio de um sistema de diálogos. Conforme apresentado no Capítulo 7, experimentos demonstraram que a folksonomia obtida a partir do aprendizado foi capaz de interpretar corretamente 72,32% dos enunciados dos diálogos num determinado domínio. Além disso, segundo (HOTH0 et al., 2006b) e (KIM et al., 2010a) as folksonomias possuem uma *semântica emergente* (STAAB et al., 2002) e (STEELS, 1998), que resulta da convergência do uso do mesmo vocabulário pelos usuários. Uma das formas de se analisar essa semântica pode ser através da exploração de tendências nas folksonomias. A detecção de tendências pode ser o monitoramento de tópicos, a sintetização de opiniões, dentre outros. Assim, com as folksonomias obtidas a partir de diálogos que esta pesquisa propõe, pode ser possível extrair tendências, como por exemplo, os assuntos abordados em um dado domínio de diálogos, num

determinado intervalo de tempo.

## 1.2. Objetivos

O objetivo principal desta pesquisa é a concepção, implementação e avaliação de um método chamado FolksDialogue, que realize o aprendizado automático de folksonomias a partir de diálogos textuais em português do Brasil.

Objetivos específicos:

- Desenvolver um protótipo computacional para avaliar o método proposto;
- Interpretar enunciados de diálogos, verificando se eles pertencem ou não ao domínio representado pela folksonomia aprendida pelo método FolksDialogue.
- Definir métricas específicas para avaliar o método proposto.

## 1.3. Hipóteses de Trabalho

As hipóteses desta pesquisa são:

*a.i) É possível construir uma folksonomia a partir de diálogos.*

*a.ii) Com a folksonomia construída, é possível indicar se enunciados de diálogos pertencem ou não ao domínio representado por ela.*

## 1.4. Contribuição Científica

A principal contribuição científica desta pesquisa é o método que realiza o aprendizado automático de folksonomias a partir de diálogos em português do Brasil. As folksonomias resultantes podem ser úteis no processo de interpretação de enunciados de diálogos, verificando se eles pertencem ou não ao domínio representado por essas estruturas. Este método, até o presente conhecimento, é o primeiro a realizar a descoberta de conhecimento a partir de diálogos textuais em português. Vale ressaltar que apesar da língua

adotada por esta pesquisa ser a portuguesa, o método proposto pode ser testado e utilizado com outras a partir de pequenas modificações, como a troca do *parser* para o de uma língua desejada.

### **1.5. Escopo**

O escopo desta pesquisa limita-se ao aprendizado automático de folksonomias a partir de diálogos textuais orientados a tarefas em português do Brasil.

### **1.6. Organização do Documento**

Este trabalho está organizado em 6 capítulos. Este é o primeiro e contém as considerações iniciais, a motivação, objetivos, hipóteses de trabalho, contribuição científica esperada e escopo.

O Capítulo 2 apresenta as definições e conceitos relacionados a folksonomias (seção 2.1) e grafos (seção 2.2), visando servir como auxílio teórico para o entendimento da presente proposta.

O Capítulo 3 tem como intuito apresentar as definições e conceitos referentes a diálogos (seção 3.1) e ao interpretador de linguagem natural (seção 3.2), servindo como fundamentação teórica para a compreensão desta proposta.

O Capítulo 4 apresenta alguns trabalhos relacionados com esta pesquisa, os quais estão organizados nas seguintes seções: Aprendizado de folksonomias (4.1), Avaliação de folksonomias (4.2) e Detecção de Tendências em folksonomias (4.3).

O Capítulo 5 apresenta o método FolksDialogue que realiza o aprendizado automático de folksonomia a partir de diálogos, e as etapas necessárias para a construção do mesmo.

O Capítulo 6 apresenta os procedimentos metodológicos utilizados para a implementação do método proposto nesta pesquisa. A implementação ocorre através de um protótipo computacional que visa validar o método FolksDialogue. Além disso, neste capítulo é descrito o corpus de diálogos utilizado como entrada do FolksDialogue, as formas de avaliação do método proposto, e também como é realizada a detecção de tendências nas folksonomias obtidas. O Capítulo 6 está organizado nas seguintes seções: Corpus de Diálogos

(6.1), Protótipo Computacional (6.2), Avaliação (6.3), Abordagem de Detecção de Tendências (6.4).

O Capítulo 7 apresenta os resultados obtidos por esta pesquisa com a realização de três experimentos com o método FolksDialogue. Na seção 7.1 é apresentado o corpus de entrada dos experimentos, e nas seções 7.2, 7.3 e 7.4 são descritos os resultados produzidos pelos experimentos realizados.

Por fim, no Capítulo 8 são apresentadas as conclusões e os trabalhos futuros.

# Capítulo 2

## Folksonomias

Neste capítulo são apresentadas as principais definições e conceitos necessários para o entendimento das folksonomias e dos grafos (forma de representação usual para as folksonomias). Entender as folksonomias é fundamental para compreender a forma de representação do conhecimento construída a partir de diálogos nesta pesquisa. Em seguida são apresentados os grafos e algumas de suas propriedades.

### 2.1. Folksonomias

Sistemas baseados em *tagging* ou rotulação colaborativa se caracterizam pela ideia de marcação de recursos ou objetos através de termos ou palavras-chave (*tags*). Tais termos são produzidos livremente pelos mais diversos usuários utilizando seus próprios vocábulos, e possuem como função servir de referência para um determinado recurso ou objeto de seus respectivos interesses. Recursos podem assumir os mais diversos tipos dependendo do sistema de rotulação que se utiliza. Exemplos de sistemas de rotulação e seus respectivos recursos são: o *Delicious*<sup>3</sup> (*URLs*), *Flickr*<sup>4</sup> (fotos) e *last.fm*<sup>5</sup> (músicas). Nesses sistemas os usuários rotulam os recursos (*URLs*, fotos ou músicas) através de uma *tag* ou rótulo, visando descrevê-los ou categorizá-los (KÖRNER et al., 2010).

Nos últimos anos houve uma crescente popularização no uso dos sistemas baseados

---

<sup>3</sup> Maiores informações em: <<http://www.delicious.com/>>

<sup>4</sup> Maiores informações em: <<http://www.flickr.com/>>

<sup>5</sup> Maiores informações em: <<http://www.last.fm/>>



em rotulação colaborativa. Os principais motivos para isso são os benefícios que esses modelos de rotulação podem oferecer aos usuários. De acordo com (GUPTA et al., 2010) alguns desses benefícios são: a recuperação futura da informação, contribuição e compartilhamento, organização de tarefas, expressão de opinião, dentre outros.

A estrutura de representação do conhecimento que emerge do processo de *tagging* ou rotulação é denominada folksonomia (XIAO et al., 2010) (CATTUTO et al., 2007) (PETERS, 2009). Segundo Thomas Van der Wal (VAN DER WAL, 2004), criador do termo “folksonomia”, a palavra deriva da junção de “*folk*” e “*taxonomia*”, ou seja, a taxonomia criada pelo povo. Uma característica importante das folksonomias é que elas são estruturas planas, ou seja, não há hierarquia (“relação entre pais e filhos”) entre suas entidades: usuários, rótulos e recursos (MATHES, 2004) (KIM et al., 2010b) (QUINTARELLI, 2005). Porém, o termo “*taxonomia*” da palavra folksonomia pode gerar confusão e dar ênfases indevidas para a característica hierárquica de uma dada classificação de termos (PETERS, 2009).

Nesta pesquisa adotamos a definição de que as folksonomias são estruturas planas. O objetivo é deixar claro a não existência de relação de subsunção entre usuários, rótulos e recursos.

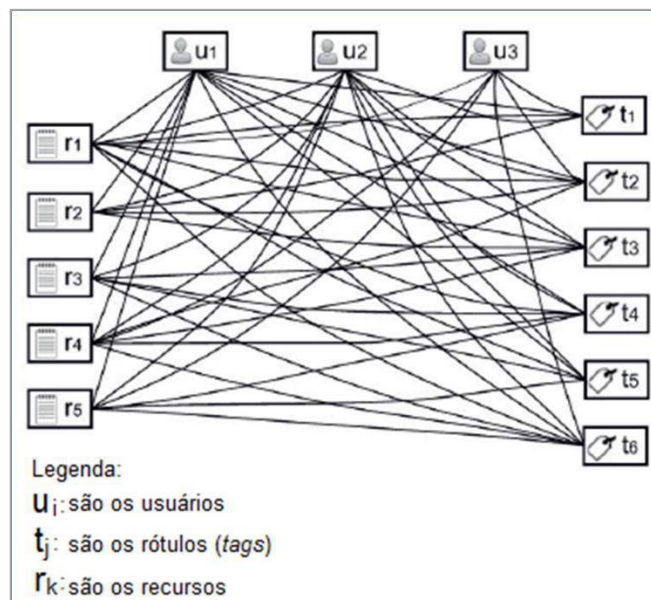
As Folksonomias podem ser definidas através de um modelo bem formalizado e aceito, denominado “modelo tripartido” (HALPIN et al., 2007) (MIKA, 2007), no qual elas são compostas por três entidades: usuários, rótulos e recursos, além de uma relação que conecta essas entidades.

Baseado na abordagem de (SCHMITZ et al., 2006), uma folksonomia pode ser definida como uma tupla  $\mathbb{F} := (U, T, R, Y)$  sendo:

- $U, T, R$  os conjuntos finitos de usuários, *tags* e recursos, respectivamente.
- $Y$  é a relação ternária entre eles, isto é,  $Y \subseteq U \times T \times R$ . Tal relação também é chamada de **atribuição**.

A *personomia*  $P_u$  de um dado usuário  $u \in U$  é a restrição em  $\mathbb{F}$  para  $u$ , ou seja,  $P_u := (T_u, R_u, I_u)$  com  $I_u := \{(t, r) \in T \times R : (u, t, r) \in Y\}$ . Em outras palavras, a personomia de um determinado usuário corresponde ao conjunto de todas as atribuições que ele gerou no processo de rotulação de um dado domínio. Com base nisso, pode-se inferir que uma folksonomia é a reunião das personomias de todos os usuários que participaram do processo de rotulação de um domínio qualquer.

A Figura 1 mostra um exemplo de folksonomia com as três entidades: usuários, rótulos e recursos, do modelo tripartido. A relação ternária  $Y$  entre as entidades é representada pelas linhas que estão conectando elas.



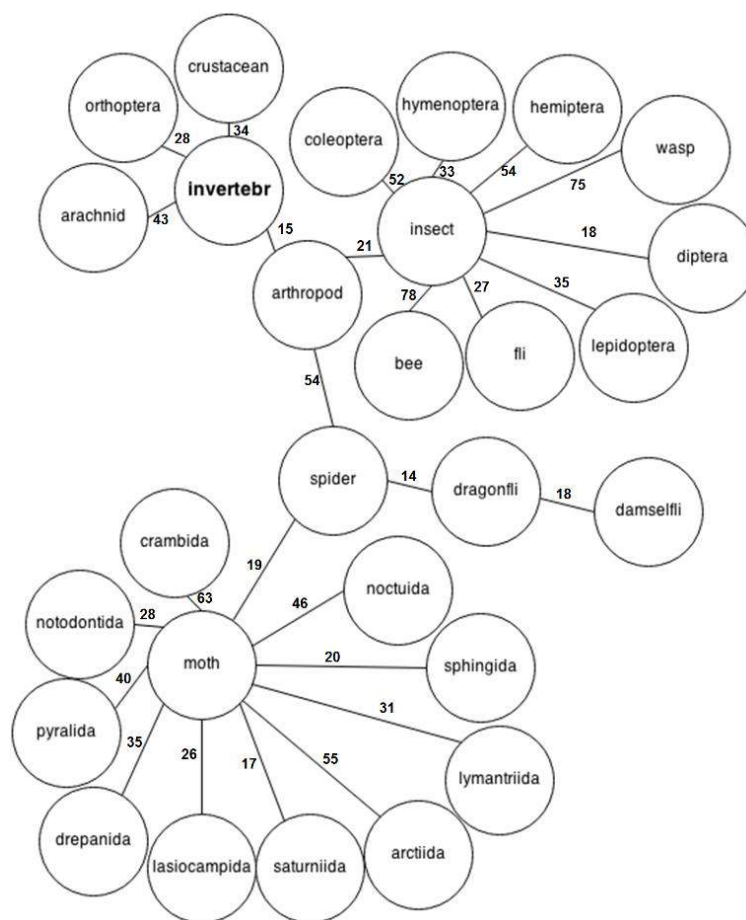
**Figura 1** - Um exemplo de folksonomia (Fonte: adaptado de DATTOLO e PITASSI (2012)).

A semântica entre conceitos está relacionada com a cognição humana. Cerca de 80% do processo cognitivo humano é formado a partir de informações visuais, do dia a dia: a coocorrência dessas informações contribui bastante para suas semânticas (WU et al., 2009). Por exemplo, a palavra “macaco” está relacionada semanticamente a “árvore”, porque os seres humanos veem frequentemente macacos vivendo em árvores. Essa coocorrência forma a relação semântica “macaco-árvore”.

Para muitos autores (BEGELMAN et al., 2006) (HALPIN et al., 2007) (SPECIA e MOTTA, 2007), o fato de duas *tags* aparecem frequentemente juntas, rotulando os mesmos recursos, é sinal da existência de um relacionamento entre elas. De acordo com (WU e ZHOU, 2011), a alta frequência de coocorrência entre *tags* não é mera coincidência, pelo contrário, ela expõe as semânticas provenientes do *tagging* colaborativo. Deste modo, é possível associar o relacionamento entre *tags* de uma dada folksonomia. Alguns autores definem nomes para os relacionamentos, como por exemplo, no contexto da construção de folksonomias a partir de códigos fontes (BEAL et al., 2014), no qual um relacionamento pode indicar se um método pertence a uma determinada classe. Outros, porém, adotam pesos nos

relacionamentos entre as *tags*, os quais são o número de recursos que elas rotulam juntas (BEGELMAN et al., 2006). Nesse caso, duas *tags* quaisquer  $t_A$ ,  $t_B$ , possuem um relacionamento  $b$  entre si, se e somente se elas apareceram juntas (rotularam os mesmos recursos) no mínimo  $x$  vezes. Sendo  $x$  o peso do relacionamento entre elas. Formalmente,  $\forall u, r, t_A, t_B \mid (u, t_A, r) \in Y \wedge (u, t_B, r) \in Y \rightarrow b(t_A, t_B) \wedge t_A \neq t_B$ .

A Figura 2 mostra um exemplo de relacionamento entre os rótulos de uma dada folksonomia pertencente ao domínio dos invertebrados. Os pesos de ligação indicam o número de recursos em que os rótulos apareceram juntos.



**Figura 2** - Exemplo de relacionamento entre rótulos de uma folksonomia (adaptado de PLANGPRASOPCHOK e LERMAN (2009)).

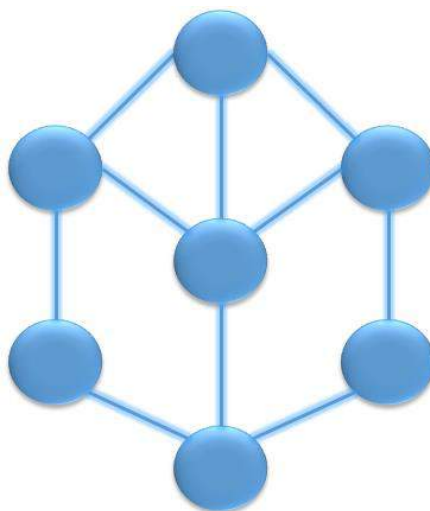
Segundo (CHOJNACKI e KLOPOTEK, 2010) as folksonomias podem ser representadas através de um grafo tripartido de usuários, *tags* e recursos. Assim, na próxima seção, serão apresentados os grafos e algumas de suas propriedades, úteis para compreender a formação das folksonomias a partir do método FolksDialogue.

## 2.2. Grafos

Folksonomias são normalmente representadas sob a forma de grafos. Segundo a definição de (SIPSER, 2007), um **grafo não-direcionado** ou simplesmente **grafo**, é um conjunto de pontos com linhas conectando tais pontos. Os pontos são denominados *vértices* ou *nós*, e as linhas são chamadas de *arestas*. Os grafos frequentemente são usados para representar dados; alguns exemplos podem ser (SIPSER, 2007):

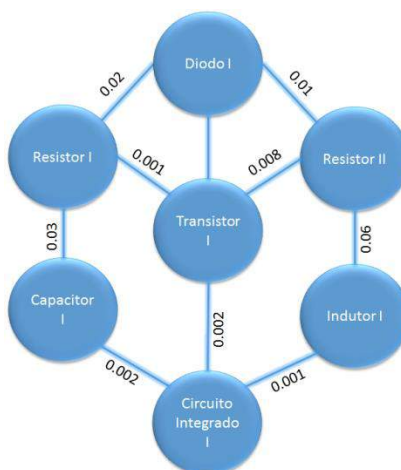
- Os nós representam cidades e as arestas, as estradas que as conectam;
- Os nós podem ser componentes elétricos, e as arestas as conexões entre eles.

Grafos podem ser visualizados através de uma representação geométrica. Nesse caso, os vértices correspondem a pontos distintos do plano em posições arbitrárias, e as arestas são linhas arbitrárias que unem esses pontos (SZWARCFITER, 1986). A Figura 3 mostra a representação geométrica com os nós (círculos) e arestas (linhas) dos grafos.



**Figura 3** - Exemplo da representação geométrica dos grafos.

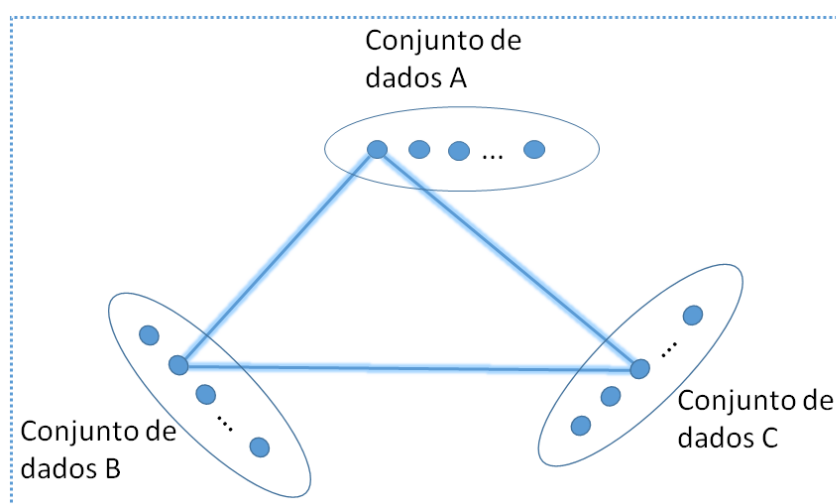
Um grafo é chamado de **grafo rotulado**, se seus nós e/ou arestas possuem algum tipo de rótulo. Um grafo rotulado é denotado por  $G := (V, E, w)$ , sendo  $V$  e  $E$  os conjuntos de vértices e arestas, respectivamente, e  $w$  o rótulo dos nós e/ou arestas. Os rótulos das arestas são também denominados **pesos**, sendo normalmente número, ou seja,  $w : E \rightarrow \mathbb{R}$ . Na Figura 4 é apresentado um exemplo de grafo rotulado, no qual seus nós representam componentes eletrônicos e as arestas são as conexões que ligam os componentes, possuindo como peso o valor da corrente elétrica que passa naquele trecho de fio.



**Figura 4** - Exemplo de grafo rotulado.

Formalmente um grafo  $G$  é denotado por:  $G := (V, E)$ , sendo que  $V$  é o conjunto de vértices e  $E$  é o conjunto de arestas de  $G$ . Dados dois vértices  $x$  e  $y$  de  $G$ , uma aresta  $e \in E$  ligando eles é denotada pelo par  $e := (x, y)$ . Neste caso, os vértices  $x$  e  $y$  são as extremidades da aresta  $e$ , sendo denominados *adjacentes*.

Chama-se de *grafo tripartido* um grafo  $G := (V, E)$ , no qual o conjunto de seus vértices  $V$  é particionado em três conjuntos de dados independentes, ou seja,  $V := A \cup B \cup C$ . A Figura 5 mostra um grafo tripartido com três tipos de vértices pertencentes aos conjuntos A, B e C.



**Figura 5** - Um grafo tripartido, com três conjuntos de vértices: A, B e C (Fonte: adaptado de CHOJNACKI e KLOPOTEK (2010)).

### 2.2.1. Representação de Folksonomias por Grafos

Uma das formas de se representar computacionalmente as folksonomias é através de um grafo tripartido de usuários, rótulos e recursos (CHOJNACKI e KLOPOTEK, 2010).

Deste modo, um grafo tripartido  $G := (V, E)$  de uma folksonomia, possuirá as seguintes características:

- O conjunto  $V$  de vértices é formado pelas três entidades usuários, rótulos e recursos, ou seja,  $V := U \cup T \cup R$ ;
- Uma aresta  $e \in E$  conecta dois nós, somente se existe uma atribuição (um usuário rotulou um recurso com uma *tag*) que relacione eles:
  - $\forall u, r E_T(u, r) \rightarrow \exists t Y(u, t, r)$  (um rótulo conectando um usuário a um recurso)
  - $\forall u, t E_R(u, t) \rightarrow \exists r Y(u, t, r)$  (um recurso conectando um usuário a um rótulo)
  - $\forall t, r E_U(t, r) \rightarrow \exists u Y(u, t, r)$  (um usuário conectando um rótulo a um recurso)

Assim,  $E := E_T \cup E_R \cup E_U$ . A Figura 6 mostra um exemplo do grafo tripartido das folksonomias. A relação ternária  $Y$  entre as entidades é representada pelas arestas (linhas) conectando elas.

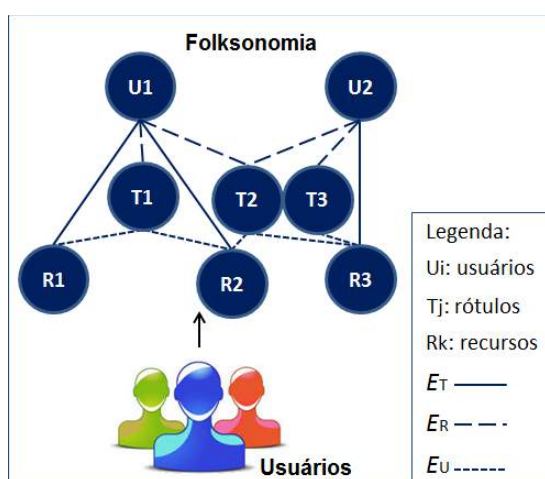
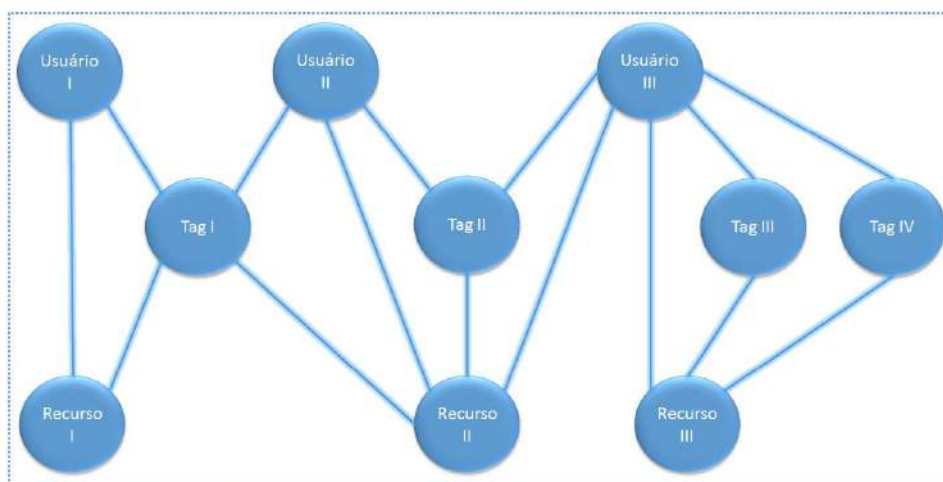


Figura 6 – Um exemplo de folksonomia.

Na Figura 7 é possível visualizar um grafo tripartido que representa parte de uma folksonomia. Algumas conclusões podem ser extraídas desse grafo:

- O *Usuário I* rotula o *Recurso I* com a *Tag I*;
- O *Usuário II* utiliza tanto a *Tag I*, quanto a *Tag II* para rotular o *Recurso II*;
- O *Usuário III* rotula o *Recurso II* com a *Tag II*, e utiliza as *Tags III* e *IV* para rotular o *Recurso III*.



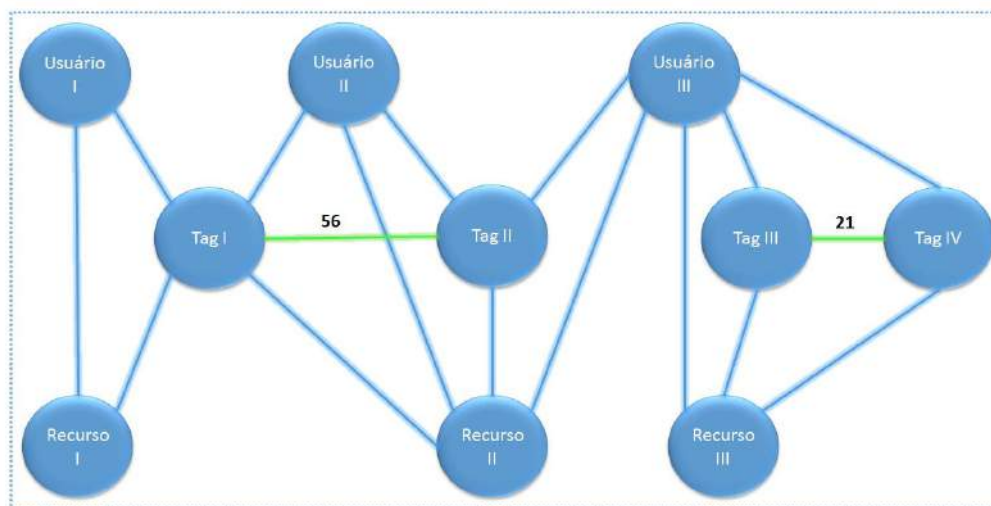
**Figura 7** - Um grafo tripartido que representa um trecho de uma folksonomia.

Em termos da utilização de relacionamentos entre os rótulos de uma folksonomia, um grafo tripartido  $G := (V, E, w_t)$  possui o componente  $w_t$ , que representa os pesos das arestas que relacionam as *tags*. Uma aresta  $e$  conecta duas *tags* quaisquer  $t_A, t_B$ , somente se elas apareceram juntas (rotularam os mesmos recursos) no mínimo  $x$  de vezes, sendo que  $e \in E$  (conjunto finito de relacionamentos) e  $x \in \mathbb{N}$ .

A Figura 8 mostra um pedaço de uma folksonomia com o relacionamento entre *tags* em destaque. Pode-se concluir que:

- As *tags I* e *II* apareceram juntas, rotulando os mesmos recursos, cinquenta e seis vezes;
- As *tags III* e *IV* rotularam juntas vinte e um recursos;
- Não existe relacionamento entre as demais *tags*, pois elas não apareceram juntas um número mínimo  $x$  de vezes.





**Figura 8** - Um trecho de uma folksonomia com os relacionamentos entre tags.

### 2.2.2. Mundo-Pequeno

A noção de *mundo-pequeno* ou também denominado *fenômeno do mundo-pequeno* surgiu num experimento proposto por (MILGRAM, 1967). Milgram propôs um experimento com o intuito de averiguar qual era a probabilidade de que duas pessoas nos Estados Unidos, selecionadas de forma aleatória, poderiam conhecer uma a outra. O autor concluiu que duas pessoas quaisquer que não se conhecem estão conectadas entre si através de seis pessoas intermediárias (em média). Um número que pode ser considerado relativamente pequeno se comparado com a população dos Estados Unidos.

Outra forma de se analisar o problema é imaginar a população de pessoas em questão como sendo uma rede social, um grafo em que cada nó representa um indivíduo e suas arestas representam as conexões ou relacionamentos entre eles (CHEN et al., 2009). Assim, o objetivo é tentar encontrar o caminho mais curto entre dois nós escolhidos de forma arbitrária.

Um *caminho*  $C$  num grafo  $G$  é uma sequência finita  $C := (v_1, v_2, v_3, \dots, v_n)$ , sendo  $n \geq 1$  e  $v_i \in V$  (o conjunto de vértices de  $G$ ).  $C$  é um caminho de  $v_1$  até  $v_n$  e seu comprimento é  $n-1$ , ou seja, o número de arestas que interligam os dois vértices.

Dados dois vértices  $v_1$  e  $v_2$  de um grafo  $G$ , se existir caminho de  $v_1$  até  $v_2$  em  $G$ , a *distância mínima*  $d$ , entre  $v_1$  e  $v_2$ , é o menor comprimento dentre todos os caminhos em  $G$  de  $v_1$  até  $v_2$ . A distância mínima entre dois nós,  $d(v_1, v_2)$ , também é denominada *caminho mais curto* entre  $v_1$  e  $v_2$ .



Recentemente houve uma associação do conceito de mundo-pequeno com diversos tipos de redes. Em (WATTS e STROGATZ, 1998), os autores mostraram que tanto redes criadas pela natureza, quanto construídas pelos homens possuem o aspecto de mundo-pequeno. Segundo (WATTS e STROGATZ, 1998) redes de energia elétrica, grafos de filmes, dentre outros, apresentam a característica de mundo-pequeno. Ainda, redes neurais de macacos e gatos, redes de sistemas de transportes e a própria internet também são exemplos de redes que possuem o fenômeno do mundo-pequeno (LATORA e MARCHIORI, 2001). Em (TACHIMORI et al., 2013), os autores verificaram que redes construídas com dados de bases médicas e práticas clínicas também possuem o fenômeno de mundo-pequeno.

Uma rede com o aspecto de mundo-pequeno pode ser definida como um grafo que tem como característica as seguintes propriedades: seu **coeficiente de agrupamento deve ser grande** e o **comprimento de caminho característico deve ser pequeno** (WATTS e STROGATZ, 1998). Tais propriedades são o que difere grafos com o fenômeno de mundo-pequeno, para grafos puramente aleatórios (conjunto de vértices, no qual as arestas estão conectadas de forma aleatória) ou puramente regulares (todos os vértices possuem o mesmo número de adjacências).

O **comprimento de caminho característico**  $Q$  é a mediana das médias das menores distâncias entre cada vértice e todos os outros, conforme representa a Equação (2.1):

$$Q = \text{mediana}\{\text{média}\{d(v1, v2)\} : v1, v2 \in V\} \quad (2.1)$$

Dado um vértice  $v$  de um grafo  $G$ , chama-se de **vizinhança de  $v$** , denotada por  $\Gamma_v$ , todos os vértices que são adjacentes a  $v$ . O **coeficiente de agrupamento** de  $v$ ,  $\gamma(v)$ , mede a proporção com que seus vizinhos estão conectados entre si. A Equação (2.2) mostra o coeficiente de agrupamento, no qual o numerador  $|E(\Gamma_v)|$  é o número de arestas que realmente existem entre os vizinhos de  $v$ , e o denominador  $(Kv(Kv - 1))/2$ , com  $Kv$  sendo o número de nós adjacentes a  $v$ , representa o total de conexões que podem existir entre os vizinhos de  $v$ .

$$\gamma(v) = \frac{|E(\Gamma_v)|}{((Kv(Kv-1))/2)} \quad (2.2)$$

O coeficiente de agrupamento de um grafo  $G$  é a média dos coeficientes de agrupamento de seus vértices, conforme Equação (2.3):

$$\gamma(G) = \text{média}\{\gamma(v)\}, \forall v \in V \quad (2.3)$$

Grafos completamente aleatórios possuem comprimento de caminho característico e coeficiente de agrupamento pequenos, enquanto que grafos puramente regulares possuem o coeficiente de agrupamento e o comprimento de caminho característico grandes. Assim, percebe-se que grafos de mundo-pequeno possuem uma topologia intermediária a grafos aleatórios e regulares, pois juntam as propriedades de ambos, como o coeficiente de agrupamento elevado e comprimento de caminho característico pequeno.

De acordo com (WATTS e STROGATZ, 1998), diz-se que o coeficiente de agrupamento de um grafo é grande, se ele for muito maior do que o coeficiente de agrupamento de um grafo aleatório, e próximo do coeficiente de agrupamento de um grafo regular (Inequação (2.4)). Já o comprimento de caminho característico, por sua vez é considerado pequeno se for muito menor do que o comprimento de caminho característico de um grafo regular, e próximo do comprimento de caminho característico de um grafo aleatório (Inequação (2.5)):

$$\gamma_{\text{Regular}} \approx \gamma_{\text{Mundo-pequeno}} \gg \gamma_{\text{Aleatório}} \quad (2.4)$$

$$Q_{\text{Aleatório}} \approx Q_{\text{Mundo-pequeno}} \gg Q_{\text{Regular}} \quad (2.5)$$

Para alguns autores (CATTUTO et al., 2007), uma característica importante dos grafos que representam as folksonomias é o fato deles possuírem o fenômeno de mundo-pequeno. Ou seja, enquanto os nós estão fortemente agrupados (coeficiente de agrupamento elevado), ainda existem conexões que ligam nós de um grupo a outros que estão longe deles (num outro agrupamento). Assim, torna-se possível caminhar por todo o grafo de modo rápido, pois o comprimento de caminho característico é pequeno. No contexto de uma folksonomia, ter um alto coeficiente de agrupamento (em relação a um grafo aleatório), pode ser um indicador de coerência no comportamento dos usuários no processo de rotulação (CATTUTO et al., 2007). Por exemplo, se um conjunto de rótulos é atribuído a certos tipos de recursos, os usuários fazem-no de modo consistente. Já no caso de uma folksonomia possuir seu comprimento de caminho característico pequeno (próximo ao de um grafo aleatório), pode ser um indicador de

que qualquer usuário, rótulo ou recurso pode ser alcançado a partir de qualquer outro nó, com apenas poucos passos em média (CATTUTO et al., 2007). Isso pode fazer com que uma busca por uma entidade (usuário, rótulo ou recurso) da folksonomia seja rápida. Deste modo, após a construção de uma folksonomia, uma das características que pode ser verificada em sua estrutura é a de mundo-pequeno.

### **2.3. Conclusão**

Neste capítulo foram abordados conceitos fundamentais que orientam este trabalho. Na seção 2.1 foi abordado sobre folksonomias. Mostrou-se o conceito onde são utilizadas e apresentou-se uma definição formal para essas estruturas de representação do conhecimento. Na seção 2.4 foram apresentados grafos. Foram descritos os principais conceitos, de que modo as folksonomias podem ser representadas através de grafos tripartidos, e também se apresentou o fenômeno de mundo-pequeno (de acordo com alguns autores, uma característica presente nas folksonomias).

# Capítulo 3

## Diálogos

Neste capítulo são apresentados os principais conceitos relacionados a diálogos e ao processo de interpretação de linguagem natural. Conhecer a estrutura e os componentes do diálogo é importante para entender a fonte dos dados utilizada para o aprendizado das folksonomias. A seguir nas seções 3.1 e 3.2 são apresentados os diálogos e o interpretador de linguagem natural, respectivamente.

### 3.1. Diálogos

O diálogo é essencialmente a interação entre falantes e ouvintes, denominados interlocutores, e é composto de atos da fala (enunciados). Segundo (AUSTIN, 1962), um enunciado não serve apenas para descrever um estado das coisas, mas também para realizar uma intenção. Ou seja, as ações realizadas por um falante através de um enunciado visam intencionalmente obter algo do ouvinte.

Todo ato de fala é ao mesmo tempo locucionário, ilocucionário e perlocucionário (AUSTIN, 1962):

- **Atos Locucionários:** Correspondem à utilização bem estruturada de uma sequência de palavras em uma determinada linguagem. São as informações passadas através dos enunciados.
- **Atos Ilocucionários:** Correspondem às ações que os falantes pretendem realizar quando produzem os enunciados, como por exemplo, pedir, cumprimentar, prometer, exigir, desculpar, censurar, entre outros.

- Atos Perlocucionários: Correspondem essencialmente às consequências ou efeito real dos atos anteriores nos ouvintes, em outras palavras, são as reações dos ouvintes em relação aos atos anteriores.

Suponha-se que alguém proferiu a frase “a sua motocicleta está trancando a passagem do meu carro”. Para esse enunciado, o ato locucionário é a informação de uma situação, o ato ilocucionário é representado pela intenção remetente ao protesto ou advertência para que a outra pessoa retire a motocicleta e por fim, o ato perlocucionário é a consequência do enunciado, de agrado ou de desagrado, podendo ou não a pessoa retirar a motocicleta.

Ainda de acordo com (AUSTIN, 1962) os atos ilocucionários possuem alguns tipos de expressões, as quais são divididas em cinco grupos:

- Expressões veridictivas: que julgam determinado assunto com base em evidências ou em razões;
  - Exemplo: Devido aos danos ao revestimento dos vasos sanguíneos, o paciente é portador de vasculite.
- Expressões exercitivas: que tomam uma decisão a favor ou contra determinada ação ou conjunto de ações.
  - Exemplo: Eu proíbo você de entrar nesta sala.
- Expressões comissivas: aquelas que comprometem o falante com o cumprimento de algo.
  - Exemplo: Eu garanto que te entregarei o carro amanhã.
- Expressões conductivas: trata-se de uma reação em relação ao destino ou conduta de outras pessoas. Podem ser agradecimentos, felicitações, saudações, entre outros.
  - Exemplo: Desejo a você um bom fim de ano.
- Expressões expositivas: sua intenção é tornar claro como a expressão do falante deve ser considerada para permanecer fiel ao seu pensamento.
  - Exemplo: A testemunha relatou que a colisão ocorreu por volta das 10h00.

Os Atos de Fala podem ser classificados em atos de fala diretos e atos de fala indiretos (SEARLE, 1969):

- Ato de Fala Direto: É definido quando realizado por meio de formas linguísticas de um determinado tipo de ato.
  - Exemplos: Que horas são? (ato de perguntar); Saia daqui (ato de ordenar); Por favor, traga-me um copo de água (ato de pedir);

- Ato de fala indireto: É definido quando realizado indiretamente por meio de formas linguísticas típicas de outro tipo de ato.
  - Exemplos: Você tem um cigarro? (pedido com aparência de pergunta), percebe-se que a pergunta não faz referência à verificação de um indivíduo ter ou não um cigarro, mas sim que uma pessoa está solicitando um cigarro a um determinado indivíduo.

Segundo (SEARLE, 1969) os atos Ilocucionários da fala podem ser classificados de acordo com seu modo de comunicação:

- Ato Ilocucionário Assertivo: Ato de fala que o falante realiza ao pronunciar um enunciado que leva alguma obrigação com os valores relativos de verdade ou falsidade. Comumente tais atos são encontrados em frases que possuem verbos assertivos e expressões verbais.
- Ato Ilocucionário Diretivo: Atos de fala a partir dos quais o falante pretende levar o ouvinte a fazer ou a dizer alguma coisa. São utilizados normalmente os verbos convidar, pedir, requerer, ordenar. Os enunciados produzidos com a intenção de levar o ouvinte a realizar algo ocorrem mais frequentemente em frases imperativas e interrogativas.
- Ato Ilocucionário Compromissivo: Ato de fala que o falante realiza com o comprometimento de realizar no futuro uma determinada ação.
- Ato Ilocucionário Expressivo: Ato de fala em que são expressos sentimentos ou emoções.
- Ato Ilocucionário Declarativo: Ato de fala que institui ou altera o estado das coisas pela simples declaração de que elas existem. Está associado a rituais, como o casamento ou o batismo. Os atos declarativos têm de obedecer às regras linguísticas específicas de uma determinada instituição e os papéis sociais do falante e do ouvinte são bem definidos.
  - Ato Ilocucionário Declarativo Assertivo: Ato Ilocucionário declarativo em que o falante tem autoridade específica, por exemplo, para excluir ou aceitar alguém num concurso ou, declarar alguém inapto para o serviço militar. Assim, estes atos de fala são simultaneamente asserções e declarações.

- Ato Ilocucionário Indireto: Ato de fala em que o falante tem a intenção de dizer algo diferente daquilo que expressa, contando com a capacidade do ouvinte em reconhecer o objetivo ilocutório do enunciado.

De acordo com (DAVIS, 1990) os atos podem também ser classificados em:

- Atos Declarativos: que transportam informação.
- Atos Interrogativos: que requerem informação.
- Atos Imperativos: que fazem uma requisição ou enviam um comando.
- Atos de Exclamação: que expressam emoção.
- Atos Performativos: que trazem à tona uma condição.

Dentre os diferentes tipos de diálogos existentes, citam-se os *diálogos orientados a tarefas*. Diálogos orientados a tarefas tem como intuito a solução de uma determinada tarefa dentro de um domínio qualquer. Tais diálogos trazem a sequência concisa de solução de uma tarefa, partindo da requisição de um solicitante a fim de realizar algo até a solução dada por outro interlocutor, o que pode ser utilizado para determinar o roteiro de solução de uma tarefa. Assim, normalmente uma das características perceptíveis nos diálogos orientados a tarefas é o fato de existirem dois tipos de interlocutores, um buscando auxílio e o outro com conhecimento do domínio, visando prestar assistência na solução de uma tarefa. Para (TRAUM e HINKELMAN, 1992), uma das características principais dos diálogos orientados a tarefas é a difusão do conhecimento, ou seja, o interlocutor com mais conhecimento *transfere a informação* (CARLETTA et al., 1997) para o que está buscando auxílio. Nesta pesquisa, em termos de adoção de nomenclatura chamam-se os interlocutores que possuem conhecimento do domínio de **atendentes** e os usuários que buscam o auxílio de **usuários**.

Segundo (IVANOVIC, 2008), existe uma diferença fundamental entre diálogos sociais ou espontâneos e diálogos orientados a tarefas. Diálogos sociais buscam manter o relacionamento entre os interlocutores, enquanto que diálogos orientados a tarefas têm como objetivo principal a realização de uma tarefa conhecida. Além disso, ainda de acordo com (IVANOVIC, 2008) conversações orientadas a tarefas tendem a ser mais estruturadas, pois focam num problema específico ou em informações específicas, ou seja, possuem um domínio bem definido. Por outro lado, diálogos sociais são menos estruturados, possuindo enunciados mais efêmeros, os quais são motivados por razões mais sutis (não possuem um domínio tão restrito).

A Figura 9 mostra um exemplo de diálogo orientado a tarefa. Nesse caso o diálogo traz

a indicação de uma solicitação de “férias” como tarefa. Esse diálogo é composto por dois interlocutores, “ATENDENTE” (detentor de conhecimento pleno do domínio) e “USUÁRIO” (quem está buscando o auxílio).

ATENDENTE:	Bom dia. Como posso ajudar?
USUÁRIO:	Gostaria de solicitar minhas férias.
ATENDENTE:	Qual o número de sua matrícula?
USUÁRIO:	12.381-4.
ATENDENTE:	Você deseja vender dias das suas férias?
USUÁRIO:	Sim.
ATENDENTE:	Quantos dias?
USUÁRIO:	10.
ATENDENTE:	Qual a data que você deseja começar a fruição?
USUÁRIO:	01/07/2014.

**Figura 9** - Exemplo de diálogo orientado a tarefa.

Para que se torne possível a interpretação computacional de enunciados é necessário o uso de um interpretador de linguagem natural, o qual será apresentado na próxima seção.

### 3.2. Interpretador de Linguagem Natural

Um dos componentes fundamentais dos sistemas de diálogos é o interpretador de linguagem natural. O papel do interpretador é analisar os enunciados de entrada dos sistemas de diálogos. A compreensão da linguagem envolve a análise semântica para determinar os significados dos constituintes (MCTEAR, 2002). A análise semântica visa definir os significados das palavras, frases, sentenças e documentos (ASSAL et al., 2010), ou extrair os significados dos enunciados (RUSSELL e NORVIG, 2003).

O processo de interpretação ocorre com base na extração dos significados dos enunciados, a qual é realizada pelas etapas do PLN mencionadas no Capítulo 1 deste trabalho. No PLN, a etapa de análise morfológica realiza a verificação da construção das palavras com suas inflexões. A análise sintática é o estudo do relacionamento entre palavras de maneira formal (JURAFSKY e MARTIN, 2008). O intuito da análise semântica é definir os significados das palavras, frases, sentenças e documentos (ASSAL et al., 2010). Segundo (RUSSELL e NORVIG, 2003) a análise semântica visa extrair os significados dos enunciados. Por fim, na análise pragmática é realizado o processamento do modo em que a linguagem é utilizada para informar, e também é analisada a forma com que os conhecimentos extraídos na



análise semântica interferem nas pessoas e em seus contextos (NEAL e SHAPIRO, 1987).

A partir disso, os significados são confrontados com uma estrutura de representação do conhecimento, como frames e ontologias. Conforme descrito no Capítulo 1, tais estruturas especificam ou representam o modelo conceitual responsável por descrever um dado domínio.

Um exemplo do processo de interpretação de enunciados de diálogos com base em frames pode ser dado por um sistema de diálogo de marcação de passagens rodoviárias. O sistema requisita ao usuário “De qual cidade você está partindo?”, a aplicação pode verificar apenas os nomes de cidades, ou talvez frases do tipo “Eu quero (sair | partir) de [nome da cidade]”. Esse sistema de marcação de passagens rodoviárias, que tem o objetivo ajudar um usuário a encontrar um horário apropriado, teria um frame com *slots* para informações sobre os horários das linhas. A Figura 10 pode corresponder ao preenchimento do frame para a entrada do enunciado “Mostre-me os horários da manhã de Londrina para Curitiba terça-feira”.

MOSTRE:
LINHAS:
ORIGEM:
CIDADE: Londrina
DATA:
DIA-DA-SEMANA: Terça-Feira
HORÁRIO:
PARTE-DO-DIA: Manhã
DESTINO:
CIDADE: Curitiba

**Figura 10** - Frame para mostra horários de linha de ônibus.

No intuito de que cada *slot* do frame seja preenchido, um *parser* extrai significado do enunciado com base em uma gramática (Figura 11), na qual as entidades semânticas são expressas.

MOSTRE	= mostre   me mostre   eu preciso
HORÁRIO	= (após   por volta   antes) HORA   manhã   tarde   noite
HORA	= uma   duas   meio dia   dezoito
ORIGEM	= de CIDADE
DESTINO	= para CIDADE
CIDADE	= Londrina   Maringá   Bandeirantes   Curitiba

**Figura 11** - Gramática para extração de significado.

### 3.3. Conclusão

Neste capítulo foram abordados conceitos fundamentais que orientam este trabalho. Na seção 3.1 foram apresentados os diálogos, os quais são compostos de atos da fala. Os atos da fala descrevem a utilização bem estruturada de uma sequência de palavras em uma determinada linguagem. Além disso, descrevem as ações que os falantes pretendem realizar quando produzem os enunciados e retratam essencialmente as consequências ou efeito real dos atos anteriores nos ouvintes. Nessa seção também foram mostrados os diálogos orientados a tarefas, os quais possuem características específicas em relação aos diálogos sociais.

Na seção 3.2 foi abordado o interpretador de linguagem natural. O interpretador é um dos componentes dos sistemas de diálogos e tem como principal função analisar semanticamente se os enunciados de entrada e saída desses sistemas pertencem ou não a um dado domínio.

# Capítulo 4

## Trabalhos Relacionados

Neste capítulo são apresentados trabalhos que se relacionam com o escopo desta pesquisa. São trabalhos que realizam o aprendizado de folksonomias, a avaliação delas no âmbito do aprendizado, e abordagens que visam a detecção de tendências em folksonomias.

Para uma melhor apresentação os trabalhos estão agrupados por temas. Estes temas são: Aprendizado de folksonomias (4.1), Avaliação de folksonomias (4.2) e Detecção de Tendências em folksonomias (seção 4.3). Além disso, com o intuito de auxiliar na visualização, na seção 4.4 é apresentada uma visão geral dos trabalhos relacionados.

Foram realizadas buscas em bases como IEEE, ACM, ScienceDirect, CiteSeerX e SpringerLink com o intuito de encontrar trabalhos com aplicação de folksonomias em diálogos. O Quadro 1 apresenta as bases pesquisadas e *strings* de busca utilizadas.

**Quadro 1** - Bases Consultadas e Strings de Busca.

<b>Bases</b>	<b>String de busca</b>
IEEE - <a href="http://ieeexplore.ieee.org/Xplore/home.jsp">http://ieeexplore.ieee.org/Xplore/home.jsp</a>	"Folksonomy" and "dialogue" "Folksonomy" and "utterance"
ACM - <a href="http://dl.acm.org/dl.cfm">http://dl.acm.org/dl.cfm</a>	"Folksonomy Network" "Building Folksonomy"
ScienceDirect - <a href="http://www.sciencedirect.com/">http://www.sciencedirect.com/</a>	"Learning Folksonomy" "Inducing Folksonomy"
CiteSeerX - <a href="http://citeseerx.ist.psu.edu/index">http://citeseerx.ist.psu.edu/index</a>	"Creating Folksonomy"
SpringerLink - <a href="http://link.springer.com/">http://link.springer.com/</a>	"Folksonomy" and "trend"

#### 4.1. Aprendizado de Folksonomias

Na literatura diversos trabalhos procuram extrair diferentes tipos de estruturas dos sistemas de *tagging* colaborativo. Dentre elas podem ser citadas: regras de associação (SCHMITZ et al., 2006), sub-bases derivadas a partir de uma base (proveniente do sistema *Delicious*) (KÖRNER et al., 2010) e grupos de *tags* relacionadas (CATTUTO et al., 2008). O principal foco desta seção é apresentar trabalhos que deixam explícito que a estrutura obtida é uma folksonomia.

Nos trabalhos de (PLANGPRASOPCHOK et al., 2010) (PLANGPRASOPCHOK et al., 2011) (PLANGPRASOPCHOK e LERMAN, 2009) os autores propõem construir folksonomias através da agregação de hierarquias pessoais (coleções no *Flickr*). A estrutura final gerada pela composição de todas as hierarquias pessoais e que representa a folksonomia é uma árvore de *tags*. Em (PLANGPRASOPCHOK et al., 2010) a técnica utilizada para a obtenção da folksonomia é baseada em agrupamento relacional com medidas de similaridade. Na abordagem de (PLANGPRASOPCHOK et al., 2011) é usado um processo de inferência probabilística em conjunto com um grupo de restrições. Já em (PLANGPRASOPCHOK e LERMAN, 2009) é adotada a estatística de coocorrência para a obtenção da folksonomia.

Implementar e avaliar três classes de algoritmos para a indução de folksonomias é o objetivo do trabalho de (STROHMAIER et al., 2012). Os três algoritmos geram estruturas hierárquicas de *tags*. O processo de obtenção das folksonomias pelos algoritmos envolve métodos de agrupamento aplicados recursivamente e o uso do algoritmo *K-Means* hierárquico. Dados de *tagging* do *BibSonomy*, *CiteULike*, *Delicious*, *Flickr*, *Last.fm* são usados nos três algoritmos.

Em (XIAO et al., 2010) os autores apresentam um modelo de correlação e coocorrência da estrutura da folksonomia visando auxiliar a recomendação em *tagging*. Nesse trabalho a folksonomia é representada por um hipergrafo de usuários, *tags* e recursos, porém os autores constroem apenas projeções dessa estrutura. Através da coocorrência entre usuários, *tags* e recursos são obtidas projeções da folksonomia, as quais são expressas através de grafos bipartidos de *tags* e usuários, recursos e usuários, e *tags* e recursos. Para a construção do modelo os autores utilizam dados de *tagging* do *Delicious* e do *MovieLens*.

No trabalho de (YOO e SUH, 2010) é proposto um método de *tagging* no qual os usuários além de escolherem as *tags* para rotular os recursos, também determinam a categoria

a que elas pertencem. Segundo os autores, o intuito de se categorizar *tags* é melhorar a semântica entre elas. Todo o processo de *tagging* realizado pelos usuários é armazenado num banco de dados. Dentro do banco as *tags* e suas respectivas categorias são armazenadas em diferentes níveis, indicando que uma é mais genérica do que a outra. Um exemplo de generalidade entre *tags* e categorias pode ser entre as palavras águia e pássaro. Um usuário utilizou no sistema como *tag* a palavra águia e como sua categoria a palavra pássaro. No método proposto isso faz de pássaro ser um termo mais genérico do que águia, ou seja, águia “é filho” de pássaro. Uma folksonomia representada através de uma hierarquia de *tags* é obtida com base no conteúdo do banco de dados do sistema.

Propor um método que encontre documentos (recursos) similares a um determinado documento (recurso) escolhido por um usuário é o objetivo de (DICHEVA E DICHEV, 2011). Nesse trabalho a folksonomia é expressa através de um hipergrafo de usuários, *tags* e recursos. No entanto, os autores constroem apenas projeções da folksonomia, derivando grafos bipartidos de documentos e *tags*, e usuários e documentos. Com base nesses grafos e a partir de um documento selecionado por um determinado usuário, a busca por documentos similares é realizada através da medida de similaridade do cosseno em combinação das *tags* e dos usuários que as rotularam. Também é usada a ordem cronológica de geração dos documentos para a análise de similaridade. Como fonte os autores adotam dados de *tagging* do *CiteULike*.

Em (LEE et al., 2009) é proposta uma forma de visualização das folksonomias, na qual se represente os relacionamentos entre as *tags*. Para a obtenção da folksonomia a ser visualizada são extraídas *tags* do *Delicious* e relacionamentos entre elas são derivados a partir do corpus da *Wikipedia*<sup>6</sup>. A folksonomia criada é representada na forma de uma hierarquia composta por relações de equivalência, similaridade e subsunção entre as *tags*.

O trabalho de (WU e ZHOU, 2011) tem como objetivo construir folksonomias e avaliar o relacionamento semântico entre suas respectivas *tags*. São construídas cinco folksonomias, as quais são expressas através de grafos de *tags*. Para a geração dos grafos são utilizadas técnicas de coocorrência em conjunto com o contexto dos recursos rotulados pelas *tags*.

Em (DATTOLO e PITASSI, 2012) a finalidade é conceber a folksonomia em termos de um sistema multiagente. Nessa abordagem a folksonomia é representada como uma

---

<sup>6</sup> Maiores informações em: <<https://www.wikipedia.org/>>

entidade semântica e dinâmica, composta de subentidades (agentes) computacionalmente autônomas. No modelo proposto cada subentidade interage com as demais através do envio de mensagens, e da reação a estímulos externos através da execução de habilidades predefinidas. Dentre as classes de subentidades definidas estão os usuários, as *tags*, os recursos e as personomias dos usuários. São utilizados dados de *tagging* do *Delicious* para a construção de um protótipo.

O sistema de *tagging* colaborativo GROUPME! (ABEL et al., 2008) além de permitir que os usuários rotulem recursos com *tags*, visa também propiciar uma funcionalidade na qual eles podem criar grupos de recursos. Para suportar esse sistema os autores propõem uma extensão do modelo tripartido (usuários, *tags* e recursos) da folksonomia, através da adição de uma quarta dimensão denominada “Grupo”. A folksonomia do GROUPME! é obtida a partir da interação dos usuários, os quais atribuem *tags* e agrupam recursos.

## 4.2. Avaliação de Folksonomias

Nesta seção são apresentadas as formas com que trabalhos que realizam o aprendizado ou a construção de folksonomias, avaliam as estruturas obtidas.

Nos trabalhos de (PLANGPRASOPCHOK et al., 2010), (PLANGPRASOPCHOK et al., 2011) (PLANGPRASOPCHOK e LERMAN, 2009) os autores comparam automaticamente a folksonomia gerada com a hierarquia do *Open Directory Project*. Em (PLANGPRASOPCHOK et al., 2010) também ocorre uma avaliação estrutural da árvore construída, balanceando profundidade e largura, e além disso, uma avaliação manual é realizada por três usuários para julgarem o caminho das hierarquias que não foram comparadas com o *Open Directory Project*. Em (PLANGPRASOPCHOK et al., 2011), os autores ainda fazem uma análise estrutural de similaridade nas hierarquias pessoais que foram agregadas para a obtenção da folksonomia, e na árvore final ocorre uma medida de integridade através da verificação do número de conflitos.

As estruturas hierárquicas geradas por (STROHMAIER et al., 2012) são avaliadas de três formas: i) comparando as hierarquias geradas pelos algoritmos com hierarquias da *Wordnet*, *Yago* e *Wikipedia*; ii) julgamento humano da qualidade das hierarquias e iii) avaliação pragmática visando analisar a qualidade de navegação da estrutura hierárquica.

No trabalho de (XIAO et al., 2010) a avaliação é realizada no cenário de

recomendação de *tags*. São selecionadas *tags* aleatoriamente e o objetivo é prever os recursos e os usuários que rotularam elas. São utilizados independentemente dois conjuntos de dados (*Delicious* e *MovieLens*), para mostrar a eficácia do método.

Em (YOO e SUH, 2010) os autores criam um protótipo para demonstrar que a folksonomia construída pode ser útil para a recomendação de documentos, porém nenhuma avaliação é realizada.

Quatro avaliações são feitas em (DICHEVA E DICHEV, 2011): i) foram testados se documentos (recursos) duplicados eram classificados como similar pelo método; ii) foi verificado se para um documento extraído de um grupo de documentos similares a ele, o algoritmo retornaria os documentos remanescentes do grupo como sendo similares; iii) uma avaliação humana indicou se os documentos retornados eram ou não similares; iv) foi analisado o impacto do comportamento dos usuários na navegação da folksonomia durante a realização de uma busca.

No trabalho de (LEE et al., 2009) ocorre uma avaliação manual elaborada por um grupo de quinze estudantes de doutorado, que avaliaram se os relacionamentos gerados entre as *tags* da folksonomia eram ou não corretos.

A partir dos cinco grafos de *tags* (folksonomias) concebidos por (WU e ZHOU, 2011) foram aplicadas algumas medidas quantitativas de similaridade, dentre elas, algumas que se baseiam em *Thesaurus* e na *Wordnet*. O intuito disso era testar o nível de relacionamento entre as *tags* de cada um dos grafos.

No modelo de concepção das folksonomias em termos de um sistema multiagente proposto por (DATTOLO e PITASSI, 2012) foi criado um protótipo com dados de *tagging* do *Delicious*, porém nenhuma avaliação foi realizada.

Em (ABEL et al., 2008) é realizada uma avaliação analisando se a quarta dimensão inserida na folksonomia (os grupos) foi aceita e passou a ser utilizada pelos usuários do sistema. Para isso, coletaram-se dados do próprio sistema de *tagging* colaborativo GROUPME! e foram constatados quantos grupos os usuários rotularam em comparação com a rotulação de recursos sozinhos.

### **4.3. Detecção de Tendências em Folksonomias**

Segundo (HOTHO et al., 2006b) (KIM et al., 2010a) as folksonomias possuem uma

*semântica emergente* (STAAB et al., 2002) (STEELS, 1998), que resulta da convergência do uso do mesmo vocabulário pelos usuários. Uma das formas de se analisar essa semântica pode ser através da exploração de tendências nas folksonomias. Alguns dos objetivos que a extração de tendências pode ter são: o monitoramento de tópicos, a sintetização de opiniões e o auxílio na navegação intuitiva da estrutura da folksonomia pelos usuários, dentre outros.

O trabalho de (HOTH0 et al., 2006b) tem como intuito descobrir tendências dentro de tópicos específicos extraídos com base em *tags*, usuários ou recursos de folksonomias. Para isso, os autores utilizam o algoritmo *FolkRank* (HOTH0 et al., 2006c), que realiza o ranqueamento de tópicos específicos a partir de *tags*, usuários ou recursos. Em seguida são comparados os ranqueamentos de folksonomias que foram geradas com dados de *tagging* (no caso do *Delicious*) de diferentes períodos de tempo. De acordo com os autores é possível descobrir ranqueamentos absolutos (por exemplo, “Quem são os *top dez*?”) e os vencedores e perdedores (“Quem cresceu/caiu mais?”). Um exemplo interessante para se entender o funcionamento da descoberta de tendências nesse trabalho, pode ser em relação ao tópico *política*. No período de um ano foram apresentadas mensalmente a evolução das *tags* que apareceram pelo menos uma vez no *top dez* desse tópico. Numa das análises feitas, os autores citam um crescimento das *tags* “*bush*” e “*election*” próximo ao dia da eleição presidencial americana de 2004.

Um dos propósitos do trabalho de (JELASSI, 2012) é inserir o tempo como uma nova dimensão na tripla (usuários, *tags* e recursos) que representa as folksonomias. Deste modo o autor propõe uma “folksonomia d” composta pela quádrupla: usuários, *tags*, recursos e datas. O intuito é poder utilizar essa nova dimensão do tempo para se analisar tendências através da verificação de quais usuários, *tags* ou recursos estão ganhando (ou perdendo) popularidade num dado intervalo de tempo. Através de dados de *tagging* de um intervalo de aproximadamente três anos, extraídos do *MovieLens* e do *Last.FM*, os autores construíram “folksonomias d”. Num dos resultados obtidos foi verificado que inicialmente o filme “Harry Potter” era visto pelos usuários como um filme para crianças, devido ao uso massivo de *tags* como “*magic*”, “*kids*” e “*witch*”. Posteriormente foi averiguado que após o lançamento de novos filmes da franquia, essas *tags* caíram em desuso e rótulos como “*fantasy*” passaram a ser utilizadas de modo elevado.



#### **4.4. Visão Geral dos Trabalhos**

Com o objetivo de auxiliar na visualização das principais diferenças entre os trabalhos relacionados, no âmbito de seus respectivos temas, construiu-se o Quadro 2. Esse quadro procura elencar os trabalhos por meio dos seguintes critérios: objetivos, o tipo da estrutura construída, a técnica adotada para o alcance dos objetivos, as fontes de dados que foram utilizadas e as avaliações realizadas. O intuito do Quadro 2 é subsidiar uma análise comparativa vinculada com a proposta desta pesquisa.

Cada trabalho apresentado no quadro está classificado de acordo com os temas descritos anteriormente neste capítulo: Aprendizado de folksonomias, Avaliação de folksonomias no contexto de seus aprendizados e Detecção de Tendências em folksonomias.

**Quadro 2** - Comparação entre os Trabalhos Relacionados.

(continua)

Tema	Abordagem	Objetivo	Tipo de estrutura construída	Técnica adotada	Fonte dos dados utilizada	Tipo de avaliação
Aprendizado/ Avaliação de folksonomias	(PLANGPRASOPCHOK et al., 2010) <i>Growing a Tree in the Forest: Constructing Folksonomies by Integrating Structured Metadata</i>	Construir uma folksonomia a partir da agregação de hierarquias pessoais dos usuários (coleções no <i>Flickr</i> e <i>bundles</i> no <i>Delicious</i> ).	Árvore de <i>tags</i> .	Agrupamento relacional que utiliza medidas de similaridade local e estrutural sobre conjuntos e coleções do <i>Flickr</i> .	Dados de <i>tagging</i> do <i>Flickr</i> .	i) Compara automaticamente a árvore gerada à hierarquia do <i>Open Directory Project</i> ; ii) É realizada uma avaliação estrutural balanceando profundidade e largura e iii) Uma avaliação manual é feita por 3 participantes para julgarem os caminhos das hierarquias que não foram comparadas em (i).
Aprendizado/ Avaliação de folksonomias	(STROHMAIER et al., 2012) <i>Evaluation of Folksonomy Induction Algorithms</i>	Implementar e avaliar três classes de algoritmos que realizam a indução de folksonomias com base em diferentes conjuntos de dados.	Nos três algoritmos implementados são geradas estruturas hierárquicas de <i>tags</i> .	Métodos de Agrupamento aplicados recursivamente, algoritmo <i>K-Means</i> hierárquico.	Dados de <i>tagging</i> do <i>BibSonomy</i> , <i>CiteULike</i> , <i>Delicious</i> , <i>Flickr</i> e <i>Last.fm</i> .	i) Avaliação comparando as hierarquias geradas pelos algoritmos com hierarquias da <i>Wordnet</i> , <i>Yago</i> e <i>Wikipedia</i> ; ii) Avaliação humana julgando a qualidade das hierarquias e iii) Propõe uma avaliação pragmática, no âmbito de analisar a qualidade de navegação da estrutura hierárquica.
Aprendizado/ Avaliação de folksonomias	(XIAO et al., 2010) <i>Towards A Correlation Cooccurrence Model Generating Approach to Folksonomy</i>	Explorar uma abordagem para gerar um modelo de correlação e coocorrência da estrutura da folksonomia com o intuito de auxiliar a recomendação em <i>tagging</i> .	As estruturas geradas são grafos bipartidos de <i>tags</i> e usuários, recursos e usuários, e <i>tags</i> e recursos.	Utiliza a coocorrência entre <i>tags</i> , usuários e recursos.	Dados de <i>tagging</i> do <i>Delicious</i> e do <i>MovieLens</i> .	A avaliação é realizada no cenário de recomendação de <i>tags</i> . São selecionadas <i>tags</i> aleatoriamente e o objetivo é prever os recursos e os usuários que rotularam elas. São utilizados dois conjuntos de dados ( <i>Delicious</i> e <i>MovieLens</i> ), para mostrar a eficácia do método.

Quadro 2 - Comparação entre os Trabalhos Relacionados.

(continuação)

Tema	Abordagem	Objetivo	Tipo de estrutura construída	Técnica adotada	Fonte dos dados utilizada	Tipo de avaliação
Aprendizado/ Avaliação de folksonomias	(YOO e SUH, 2010) <i>User-Categorized Tags to Build a Structured Folksonomy</i>	Propor um método de <i>tagging</i> no qual além dos usuários escolherem a <i>tag</i> , também escolhem a categoria que ela pertence. O intuito é melhorar a semântica entre <i>tags</i> . Através desse método, uma folksonomia é construída com os dados de <i>tagging</i> gerados pelos usuários.	É criada uma estrutura hierárquica entre <i>tags</i> .	O processo de <i>tagging</i> realizado pelos usuários é armazenado num banco de dados. Com base nas <i>tags</i> e em suas categorias são definidos diferentes níveis no banco de dados, indicando que uma categoria é mais genérica do que uma <i>tag</i> . A partir do conteúdo do banco de dados, a folksonomia é construída na forma de uma hierarquia de <i>tags</i> .	São usados dados de <i>tagging</i> que deveriam ser provenientes dos usuários que utilizassem o protótipo proposto. No entanto, não são especificados se foram obtidos dados reais de usuários ou outras fontes.	Cria um protótipo dizendo que a folksonomia construída com o método proposto pode ser útil para a recuperação de documentos, porém nenhuma avaliação é realizada.
Aprendizado/ Avaliação de folksonomias	(PLANGPRASOPCHOK et al., 2011) <i>A Probabilistic Approach for Learning Folksonomies from Structured Data</i>	Construir uma folksonomia a partir da agregação de hierarquias pessoais dos usuários (coleções no <i>Flickr</i> , <i>bundles</i> no <i>Delicious</i> ).	Árvore de <i>tags</i> .	Processo de inferência probabilística que incorpora informações estruturais (hierarquias pessoais dos usuários) através de restrições, evitando deste modo loop na estrutura e forçando com que ela seja uma hierarquia.	Dados de <i>tagging</i> do <i>Flickr</i> .	i) Compara automaticamente a árvore gerada à hierarquia do <i>Open Directory Project</i> ; ii) É realizada uma avaliação estrutural analisando a similaridade entre as pequenas hierarquias que formam a árvore final. Também na avaliação estrutural são verificados o número de conflitos na estrutura, visando medir a integridade da árvore.

Quadro 2 - Comparação entre os Trabalhos Relacionados.

(continuação)

Tema	Abordagem	Objetivo	Tipo de estrutura construída	Técnica adotada	Fonte dos dados utilizada	Tipo de avaliação
Aprendizado/ Avaliação de folksonomias	(PLANGPRASOPCHOK e LERMAN, 2009) <i>Constructing Folksonomies from User-Specified Relations on Flickr</i>	Construir uma folksonomia a partir da agregação de hierarquias pessoais dos usuários (coleções no <i>Flickr</i> , <i>bundles</i> no <i>Delicious</i> ).	Árvore de <i>tags</i> .	Abordagem estatística de coocorrência para agregar as hierarquias pessoais dos usuários numa árvore final.	Dados de <i>tagging</i> do <i>Flickr</i> .	Compara automaticamente a árvore gerada à hierarquia do <i>Open Directory Project</i> .
Aprendizado/ Avaliação de folksonomias	(DICHEVA E DICHEV, 2011) <i>Can Collective Use Help for Searching?</i>	Criar um método que encontre documentos (recursos) similares a um determinado documento (recurso) escolhido por um usuário.	A partir de um hipergrafo que representa a estrutura da folksonomia, são derivados grafos bipartidos de documentos e <i>tags</i> e, usuários e documentos.	A similaridade entre dois documentos é calculada pela medida de similaridade do cosseno, combinando as <i>tags</i> e os usuários que as rotularam. Também é utilizada a recentidade dos documentos na análise de similaridade.	Dados de <i>tagging</i> do <i>CiteULike</i> .	i) Foram testados se documentos duplicados eram classificados como similar pelo método; ii) Dado um documento foi verificado se o algoritmo retornaria como similares, os documentos pertencentes a um mesmo grupo iii) Avaliação humana indicando se os documentos retornados são ou não similares; iv) Foi analisado o impacto do comportamento dos usuários na navegação da folksonomia durante a realização de uma busca.
Aprendizado/ Avaliação de folksonomias	(LEE et al., 2009) <i>FolksoViz: A Semantic Relation-Based Folksonomy Visualization Using the Wikipedia Corpus</i>	Gerar uma forma de visualização das folksonomias que represente os relacionamentos entre <i>tags</i> .	Hierarquia de <i>tags</i> composta por relações de equivalência, similaridade e subsunção.	O método extrai <i>tags</i> do <i>Delicious</i> e cria relações entre elas com base no corpus da <i>Wikipedia</i> . São criadas três tipos de relações: equivalência, similaridade e subsunção.	Dados de <i>tagging</i> do <i>Delicious</i> e corpus da <i>Wikipedia</i> .	Avaliação manual realizada por um grupo de quinze estudantes de doutorado, que avaliaram se os relacionamentos gerados entre as <i>tags</i> da folksonomia eram ou não corretos.

Quadro 2 - Comparação entre os Trabalhos Relacionados.

(continuação)

Tema	Abordagem	Objetivo	Tipo de estrutura construída	Técnica adotada	Fonte dos dados utilizada	Tipo de avaliação
Aprendizado/ Avaliação de folksonomias	(WU e ZHOU, 2011) <i>TAGS ARE RELATED: MEASUREMENT OF SEMANTIC RELATEDNESS BASED ON FOLKSONOMY NETWORK</i>	Construir folksonomias e avaliar o relacionamento semântico entre as respectivas <i>tags</i> delas.	Grafos de <i>tags</i> .	São construídos cinco grafos de <i>tags</i> através de técnicas de coocorrência e com base nos contextos dos recursos que elas rotulam.	Dados de <i>tagging</i> do <i>Delicious</i> .	A partir dos cinco grafos gerados foram aplicadas medidas quantitativas de similaridade, algumas baseadas em <i>Thesaurus</i> e na <i>Wordnet</i> , visando testar o nível de relacionamento entre as <i>tags</i> de cada um dos grafos.
Aprendizado/ Avaliação de folksonomias	(DATTOLO e PITASSI, 2012) <i>Folkview: A Multi-agent System Approach to Modeling Folksonomies</i>	Conceber a folksonomia em termos de um sistema multiagente. O intuito é prover uma estrutura dinâmica, provendo por exemplo, a habilidade dos usuários alterarem em nível global as <i>tags</i> de suas respectivas personomias.	A folksonomia é concebida como uma entidade semântica e dinâmica. Essa estrutura é organizada como um universo de subentidades (agentes) computacionalmente autônomas.	Cada subentidade da folksonomia interage com as demais através do envio de mensagens e da reação à estímulos. Foram definidas oito classes de subentidades, dentre elas, <i>tags</i> , usuários, recursos e personomia.	Dados de <i>tagging</i> do <i>Delicious</i> .	Os autores construíram um protótipo com dados de <i>tagging</i> do <i>Delicious</i> , porém não existe avaliação da proposta.
Aprendizado/ Avaliação de folksonomias	(ABEL et al., 2008) <i>A NOVEL APPROACH TO SOCIAL TAGGING: GROUPME!</i>	Apresentar um sistema de <i>tagging</i> que além de permitir que usuários rotulem recursos com <i>tags</i> , também propicie que eles criem grupos de recursos.	A estrutura derivada pelo sistema proposto é uma folksonomia (usuários, <i>tags</i> e recursos) estendida com uma quarta dimensão “grupos”.	Folksonomia com quatro dimensões gerada pela interação dos usuários com o sistema, através da rotulação de recursos e da criação de grupos. Os autores não deixam explícito qual estrutura computacional armazena a folksonomia.	Dados de <i>tagging</i> gerados pelos usuários que utilizam do sistema apresentado, “ <i>GROUPME!</i> ”.	Foi analisado se a dimensão inserida na folksonomia (os grupos) foi aceita e passou a ser utilizada pelos usuários do sistema. Para isso, coletaram-se dados do próprio <i>GROUPME!</i> e foram constatados quantos grupos os usuários rotularam em comparação com a rotulação de recursos sozinhos.

Quadro 2 - Comparação entre os Trabalhos Relacionados.

(continuação)

Tema	Abordagem	Objetivo	Tipo de estrutura construída	Técnica adotada	Fonte dos dados utilizada	Tipo de avaliação
Detecção de Tendências em folksonomias	(HOTH0 et al., 2006b) <i>Trend Detection in Folksonomies</i>	Descobrir tendências dentro de tópicos específicos extraídos com base em <i>tags</i> , usuários ou recursos de folksonomias.	Lista de ranques.	Com base no algoritmo <i>FolkRank</i> realiza o ranqueamento de tópicos específicos a partir de <i>tags</i> , usuários ou recursos. Após isso, compara os ranques de folksonomias obtidas em diferentes intervalos de tempo e analisa as mudanças nas tendências.	Dados de <i>tagging</i> do <i>Delicious</i> .	É realizado um experimento no qual a partir de dados de <i>tagging</i> obtidos em diferentes períodos de tempo do sistema <i>Delicious</i> são construídos ranques. Após isso, são feitas análises visuais comparando as mudanças dos dados dos ranques.
Detecção de Tendências em folksonomias	(HOTH0 et al., 2006b) <i>Trend Detection in Folksonomies</i>	Descobrir tendências dentro de tópicos específicos extraídos com base em <i>tags</i> , usuários ou recursos de folksonomias.	Lista de ranques.	Com base no algoritmo <i>FolkRank</i> realiza o ranqueamento de tópicos específicos a partir de <i>tags</i> , usuários ou recursos. Após isso, compara os ranques de folksonomias obtidas em diferentes intervalos de tempo e analisa as mudanças nas tendências.	Dados de <i>tagging</i> do <i>Delicious</i> .	É realizado um experimento no qual a partir de dados de <i>tagging</i> obtidos em diferentes períodos de tempo do sistema <i>Delicious</i> são construídos ranques. Após isso, são feitas análises visuais comparando as mudanças dos dados dos ranques.

Quadro 2 - Comparação entre os Trabalhos Relacionados.

(conclusão)

Tema	Abordagem	Objetivo	Tipo de estrutura construída	Técnica adotada	Fonte dos dados utilizada	Tipo de avaliação
Detecção de Tendências em folksonomias	(JELASSI, 2012) <i>A quadratic approach for trend detection in folksonomies</i>	Inserir o “tempo” como uma nova dimensão na tripla das folksonomias e com isso realizar a análise de tendências num dado intervalo de tempo.	“Folksonomias d”, as quais são folksonomias que possuem uma quarta dimensão, no caso o “tempo”.	Com base na dimensão temporal das “folksonomias d”, verificam-se quais usuários, tags ou recursos ganharam (ou perderam) popularidade no período representado por essa nova dimensão.	Dados de <i>tagging</i> do <i>MovieLens</i> e do <i>Last.fm</i> .	Foi feito um experimento com dados do <i>MovieLens</i> e do <i>Last.fm</i> . O intuito era analisar visualmente o ganho ou perda de popularidade das tags usadas para rotular determinados filmes. Também foi realizada uma análise de popularidade de usuários e seus amigos ao longo do tempo, dentro de um determinado gênero musical.

Com a organização dos trabalhos relacionados no formato do Quadro 2 foi possível inferir algumas conclusões relacionadas a cada coluna dele. Na seção a seguir são detalhadas tais conclusões.

#### **4.5. Análise dos Trabalhos Relacionados e Conclusão**

Neste capítulo foram apresentados trabalhos da literatura que estão relacionados com os temas do escopo desta pesquisa.

Como critérios de seleção e avaliação dos trabalhos foram utilizados seus objetivos, o tipo da estrutura construída, a técnica adotada para o alcance dos objetivos, as fontes de dados que foram utilizadas, e as avaliações realizadas.

Com base nesses critérios é possível realizar uma análise crítica dos trabalhos apresentados. Em relação ao aprendizado/construção das folksonomias, a grande parte dos trabalhos relacionados deriva estruturas na forma de árvores ou hierarquias de *tags*. O fato de alguns autores interpretarem as folksonomias como uma estrutura hierárquica, pode ser atribuído ao termo “taxonomia” da palavra folksonomia, o qual pode gerar confusão e dar ênfases indevidas para a característica hierárquica de uma dada classificação de termos (PETERS, 2009), conforme descrito na seção 2.1 do Capítulo 2 deste trabalho. Existem trabalhos ainda que geram grafos bipartidos, os quais são compostos apenas por duas dimensões da tripla: usuários, *tags* e recursos. Além disso, todos os trabalhos utilizam para a construção de suas folksonomias fontes de dados provenientes de sistemas de *tagging* colaborativo, como o *Delicious* e o *CiteULike*. No entanto, até o presente momento não foi encontrada nenhuma abordagem que aproveite a dimensão social que as folksonomias e os diálogos possuem em comum, e realize o aprendizado dessas estruturas a partir de tais interlocuções. Deste modo, atualmente não existem trabalhos que utilizem as características dos diálogos orientados a tarefas para a obtenção de folksonomias, como por exemplo, o fato de haverem dois tipos de interlocutores, um buscando auxílio e o outro visando prestar assistência na solução de uma determinada tarefa.

A maioria dos trabalhos que avaliam as folksonomias no âmbito de seu aprendizado, comparam as estruturas obtidas com padrões ouro, tais como *Open Directory Project* e *Wordnet*. Existem também avaliações baseadas nas opiniões de avaliadores humanos, que verificam se os relacionamentos gerados pelos métodos são ou não consistentes. O fato é que



nenhuma dessas abordagens avalia as folksonomias como uma representação do conhecimento de um dado domínio, que interpreta se enunciados de diálogos pertencem ou não a tal domínio.

Os trabalhos de detecção de tendências em folksonomias adotam abordagens distintas. Existe o uso de um algoritmo de ranqueamento e também a inserção de uma nova dimensão (no caso o tempo) na tripla das folksonomias. O que diferencia a abordagem de detecção de tendências desses trabalhos da literatura, para o que é descrito nesta pesquisa é que os trabalhos da literatura não utilizam características ou o conteúdo dos recursos nas detecções de tendências. No caso de recursos como diálogos, características específicas podem ser utilizadas a fim de se auxiliar na descoberta de tendências nas folksonomias. Detalhes no Capítulo 6.

## Capítulo 5

# Um Método para o Aprendizado Automático de Folksonomias a partir de Diálogos

Este capítulo apresenta o método FolksDialogue para o aprendizado automático de folksonomias a partir de diálogo orientado a tarefas em português do Brasil.

Nesta pesquisa propõe-se representar o Modelo Conceitual através das folksonomias, as quais são estruturas de representação do conhecimento simples de se implementar e utilizar (ECHARTE et al., 2004). Folksonomias são estruturas que emergem do processo de rotulação colaborativo, o qual corresponde às atribuições de rótulos a recursos por usuários. Um dos benefícios do processo de rotulação é que os usuários não necessitam ter experiência ou habilidades específicas (HOTHO et al., 2006a). Além disso, o vocabulário que compõe as folksonomias reflete diretamente a linguagem dos usuários comuns que participam do processo (QUINTARELLI, 2005). Uma das características interessantes das folksonomias é a sua dimensão social (usuários), a qual também é um aspecto presente em diálogos, resultante do processo de interação entre os seres humanos. Neste trabalho é proposta a obtenção de folksonomias a partir de diálogos, dessa maneira, as entidades usuários, rótulos e recursos passam a ser representadas do seguinte modo: os usuários correspondem aos interlocutores, os rótulos são termos extraídos dos enunciados e os recursos são os próprios enunciados dos diálogos.

Em comparação com os métodos citados no Capítulo 4, a relevância do método proposto nesta pesquisa para o aprendizado de folksonomias a partir de diálogos, é que ele estende o modelo tripartido das folksonomias utilizando características dos diálogos orientados a tarefas. Deste modo, o método se torna único e específico para a realização do

aprendizado a partir de diálogos.

A seguir na seção 5.1 será apresentada uma extensão da definição formal de folksonomias do Capítulo 2, no âmbito da obtenção delas a partir de diálogos orientados a tarefas. Posteriormente se apresentará o FolksDialogue: um método para o aprendizado automático de folksonomias a partir de diálogo orientado a tarefas em português do Brasil, e na sequência as etapas para a sua obtenção (seção 5.2).

### **5.1. Definição Formal de Folksonomias obtidas a partir de Diálogos**

Esta seção propomos uma extensão da definição formal do modelo tripartido das folksonomias apresentada por (SCHMITZ et al., 2006), a qual foi mostrada no Capítulo 2. O intuito de se estender essa definição é torná-la compatível para a obtenção de folksonomias a partir de diálogos orientados a tarefas. Vale ressaltar, conforme explicitado no Capítulo 2, que uma das características de diálogos orientados a tarefas é a existência de dois tipos de interlocutores, um buscando auxílio e o outro com conhecimento do domínio, visando prestar assistência na solução de uma tarefa.

Nesta pesquisa, as entidades usuários, rótulos e recursos das folksonomias são representados do seguinte modo: **os usuários correspondem aos atendentes** dos diálogos orientados a tarefas, **os rótulos são os substantivos dos enunciados gerados pelos atendentes**, e por fim, **os recursos são os próprios enunciados dos atendentes**. O fato de serem utilizados os atendentes, os substantivos de seus enunciados, e seus próprios enunciados como usuários, rótulos e recursos das folksonomias, respectivamente, é devido à suposição de que os atendentes são os interlocutores que possuem conhecimento pleno de um dado domínio. Em oposição, os interlocutores do tipo usuário buscam auxílio para a solução de uma tarefa, podendo requisitar qualquer coisa, mesmo que ela esteja fora do domínio compreendido pelos atendentes. Assim, a folksonomia estaria divergente da representação de determinado domínio, o que poderia gerar um problema para interpretação de enunciados de diálogos, no âmbito de verificar se eles pertencem ou não a tal domínio.

O processo de rotulação colaborativa se caracteriza pela atribuição de rótulos a recursos ou objetos. De acordo com (EMBLEY e THALHEIM, 2011), para ser possível a distinção entre objetos, a linguagem natural atribui nomes a eles, denominando-os de *substantivos*. Assim, justifica-se apenas a utilização de substantivos dos enunciados dos

atendentes como rótulos dos recursos das folksonomias. Além disso, em sistemas de rotulação colaborativos os usuários normalmente utilizam substantivos, os quais servem de representação para objetos como *casa*, *avião* e *violino*. Segundo (SPITERI, 2007), no sistema *Delicious*, os objetos representam a grande maioria das rotulações realizadas pelos usuários chegando a 76% do total. Já a predominância dos substantivos usados como classe gramatical de rotulação é ainda maior, com um domínio de 88% do sistema (ANDREWS et al., 2010) e (SPITERI, 2007).

A seguir são apresentadas as definições necessárias para a compreensão de uma folksonomia obtida a partir de diálogos orientados a tarefas, a qual é proposta por esta pesquisa.

**Definição 5.1.** Um **subconjunto de usuários**  $l$  pertence a um dado atendente  $a$ , e é composto por todos os usuários com quem ele dialogou num dado domínio. Vale ressaltar que um atendente para esta pesquisa é o interlocutor do diálogo orientado a tarefas que tem o conhecimento pleno de um dado domínio, conforme descrito na seção 3.1 deste documento. Cada atendente possui um e somente um subconjunto de usuários. Formalmente, seja:

- $A$  o conjunto finito de todos os atendentes (sendo  $a$  um atendente pertencente a  $A$ );
- $U$  o conjunto finito de todos os usuários (sendo  $u$  um usuário pertencente a  $U$ );
- $D$  um dado corpus de diálogos (sendo  $d$  um diálogo pertencente a  $D$ );
- $Du$  é uma função  $Du: A \times D \rightarrow U$  que retorna o usuário atendido por um atendente num dado diálogo.
- $Enu$  o conjunto de enunciados de todos os diálogos.

O subconjunto de usuários de um atendente  $a$  (sendo  $a$  uma constante) pode ser definido pelo predicado:  $l: \forall d((a, d) \in A \times D) \rightarrow l(Du(a, d))$ .

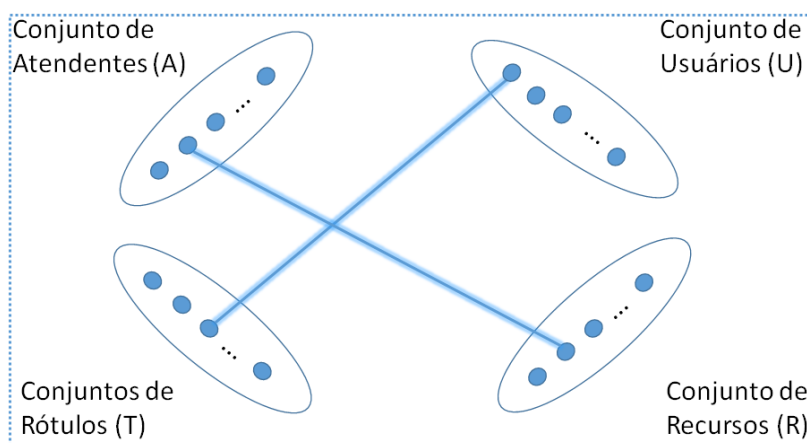
**Definição 5.2.** Formalmente uma folksonomia obtida a partir de diálogos orientados a tarefas pode ser definida como uma tupla  $\mathbb{F} := (A, T, R, U, Y')$ , sendo:

- $A$  o conjunto finito de usuários da folksonomia, ou seja, os atendentes dos diálogos (usuários que possuem conhecimento pleno do domínio);
- $T$  o conjunto finito de rótulos ou *tags*, ou seja, os substantivos dos enunciados que os atendentes produziram ao longo dos diálogos;

- $R$  o conjunto finito de recursos da folksonomia, ou seja, correspondem aos enunciados dos atendentes;
- $U$  o conjunto finito dos usuários;
- $Y'$  é a relação quaternária entre eles, isto é,  $Y' \subseteq A \times T \times R \times U$ . Tal relação também é chamada de **atribuição**.

Deste modo, percebe-se que a folksonomia obtida a partir de diálogos orientados a tarefas é representada por um *modelo quadripartido* (Figura 12), pois possui quatro dimensões (atendentes, rótulos, recursos, e subconjuntos de usuários) em comparação com o *modelo tripartido* (usuários, rótulos, e recursos) usado tradicionalmente usado para representar folksonomias (Capítulo 2). Isto é o resultado da utilização dos *usuários* ( $u \in U$ ), os quais formam os subconjuntos  $l$  e são uma característica dos diálogos orientados a tarefas. Tais *usuários* são os interlocutores do tipo *usuário* desses diálogos, e visam buscar auxílio nas soluções de tarefas (descrito na seção 3.1 deste trabalho). Além das quatro dimensões  $A$ ,  $T$ ,  $R$ ,  $U$ , a tupla  $F$  da folksonomia proposta nesta pesquisa, é composta também pela relação de atribuição  $Y'$ , a qual realiza as devidas conexões entre tais conjuntos.

A *personomia*  $P_a$  de um dado atendente  $a \in A$  é a restrição em  $F$  para  $a$ , ou seja,  $P_a := (T_a, R_a, l_a, I_a)$  com  $I_a := \forall t, r, u I_a(t, r, u) \rightarrow (a, t, r, u) \in Y'$ . Em outras palavras, a personomia de um determinado atendente corresponde ao conjunto de todas as atribuições obtidas a partir dos enunciados produzidos por ele num dado domínio. Com base nisso, deduz-se que uma folksonomia construída com diálogos orientados a tarefas é a reunião das personomias de todos os atendentes que participaram do processo de interlocução num domínio qualquer.



**Figura 12** – Modelo quadripartido que representa a folksonomia obtida a partir de diálogos.

Além disso, para a folksonomia aprendida a partir de diálogos, a qual é proposta nesta pesquisa, adota-se a noção de relacionamento entre seus rótulos (descrita no Capítulo 2 deste trabalho). Deste modo, dois rótulos quaisquer  $t_A, t_B$  da folksonomia possuirão um relacionamento  $b \in B$  (conjunto de relacionamento entre rótulos) entre eles, se e somente se, estes rótulos aparecerem juntos (rotularem os mesmos recursos) um número mínimo  $x$  de vezes. Formalmente,  $\forall a, u, r, t_A, t_B b(t_A, t_B) \rightarrow ((a, t_A, r, u) \in Y' \wedge (a, t_B, r, u) \in Y' \wedge t_A \neq t_B \wedge w(t_A, t_B) \geq x)$ .

Já o peso  $w$  (descrito no Capítulo 2), adotado no âmbito desta pesquisa, do relacionamento entre dois dados rótulos, é o **número de diálogos em que eles aparecem juntos**. Sendo que para dois rótulos serem considerados como “aparecendo juntos” eles não precisam estar num mesmo enunciado de um dado diálogo, eles podem pertencer a enunciados distintos, porém devem pertencer ao mesmo diálogo. Formalmente, o peso  $w$  do relacionamento entre dois rótulos  $t_A, t_B$  pode ser definido pela função:  $w: T \times T \rightarrow \mathbb{N}$ , sendo  $\mathbb{N}$  o conjunto dos naturais.

Após definir formalmente uma folksonomia obtida a partir de diálogos orientados a tarefas, o próximo passo é a apresentação do FolksDialogue, um método para o aprendizado automático de folksonomias a partir de diálogo orientado a tarefas em português do Brasil. Na sequência a seção 5.2 apresenta o método proposto por esta pesquisa.

## 5.2. O Método FolksDialogue

O método descrito nesta seção foi elaborado de acordo com os seguintes pressupostos:

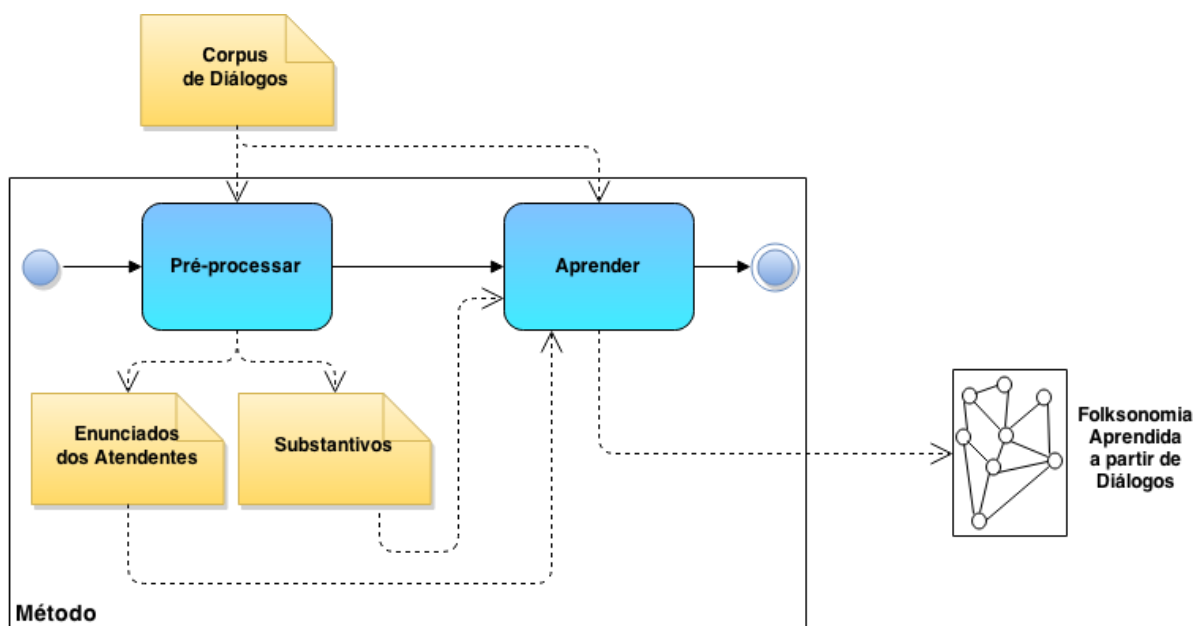
- O Modelo Conceitual deve ser representado através de folksonomias;
- A construção da folksonomia é realizada de forma automática;
- Os diálogos textuais orientados a tarefas em português do Brasil são fonte de conhecimento para a construção do modelo;
- Adota-se que cada diálogo orientado à tarefa, é formado por um interlocutor distinto do tipo *usuário*. No entanto, um mesmo interlocutor do tipo atendente pode aparecer em diferentes diálogos;
- São utilizados apenas substantivos como rótulos (conforme seção 5.1);
- Para o aprendizado das folksonomias, somente são adotados os enunciados dos

atendentes (conforme seção 5.1).

Vale destacar que apesar do método ter sido testado apenas em português do Brasil (Capítulo 7), nada impede o uso de outras linguagens a partir de pequenas modificações, como a troca do *parser* para o da língua desejada.

A Figura 13 apresenta o diagrama de atividades do método. De acordo com a figura, verifica-se que o método *FolksDialogue* é composto por duas atividades: *Pré-processar* e *Aprender*. Cada uma das atividades é composta por suas próprias **etapas** (atividades), as quais são mostradas na Figura 14 (atividade de Pré-processar) e na Figura 15 (atividade de Aprender).

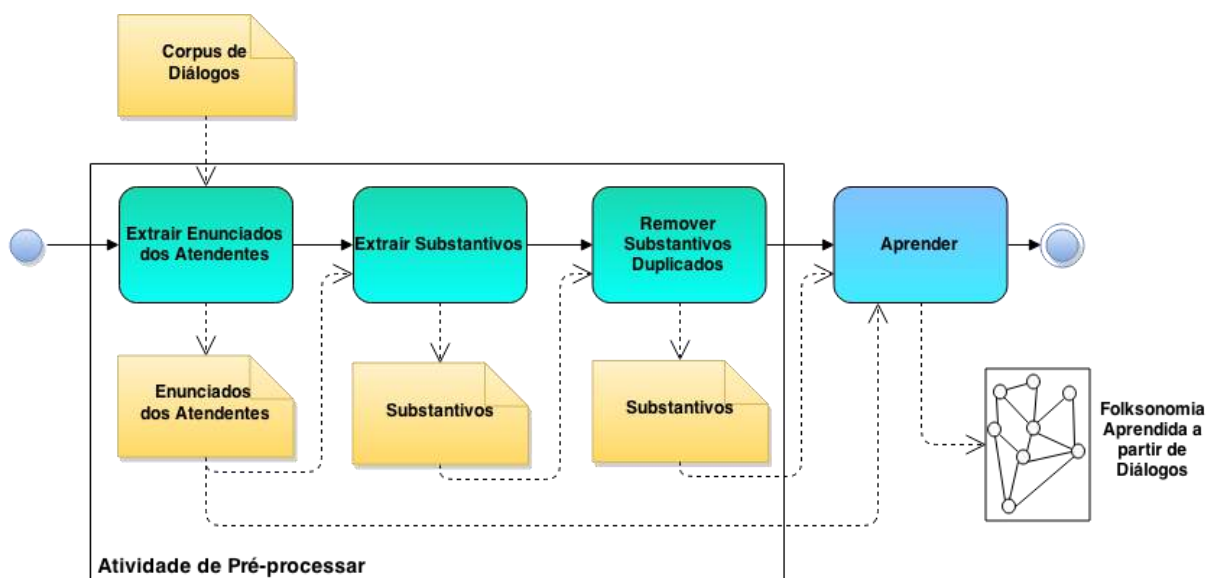
Porém, antes da execução do método existe a necessidade da entrada de um corpus de diálogos de um domínio qualquer. Para melhorar o entendimento do método proposto e das etapas que o compõem, no Quadro 3 são mostrados exemplos de diálogos orientados a tarefas, os quais serão usados como corpus de entrada para exemplificar a execução do *FolksDialogue*.



**Figura 13** - Diagrama de atividades representando o fluxo do método.

No Quadro 3, a coluna **ID do Diálogo** é um número de identificação único que representa cada um dos diálogos do corpus. Cada diálogo é formado por um conjunto de enunciados dispostos em sequência temporal, iniciando pelo enunciado mais antigo. Os

enunciados de um dado diálogo estão presentes no campo **Enunciado**. Cada enunciado está identificado de acordo com o interlocutor que o gerou (campo **Tipo de Interlocutor**), sendo: “ $a_i$ ” (tal que  $i \in \mathbb{N}^*$ ) a identificação para interlocutores do tipo atendente e “ $u_j$ ” (tal que  $j \in \mathbb{N}^*$ ) a identificação para interlocutores do tipo *usuário*.



**Figura 14** – Etapas da atividade de Pré-processar que compõe o método proposto.

Após a entrada do corpus de diálogos, inicia-se o método FolksDialogue. A atividade de Pré-processar tem como intuito receber o corpus de entrada de diálogos e tornar ele apto para ser utilizado no restante do processo. De acordo com a Figura 14, as etapas que compõem a atividade de Pré-processar são: Extrair Enunciados dos Atendentes, Extrair Substantivos e Remover Substantivos Duplicados.

**Quadro 3** - Exemplos de diálogos orientados a tarefas.

(continua)

ID do Diálogo	Tipo de Interlocutor	Enunciado
1	$u_1$	Boa tarde.
1	$a_1$	Boa tarde.
1	$u_1$	Eu quero fazer uma pergunta.
1	$a_1$	Pois não, pode fazer.
1	$u_1$	Se eu quiser tirar férias, preciso enviar alguma solicitação por escrito?
1	$a_1$	Precisa enviar uma cópia do documento intitulado Planilha de Férias indicando se você deseja converter 10 dias de férias em pecúnia ou não e indicando o início do período de fruição com a sua assinatura e a assinatura da sua chefia.



Quadro 3 - Exemplos de diálogos orientados a tarefas.

(conclusão)

ID do Diálogo	Tipo de Interlocutor	Enunciado
1	$a_1$	Esse documento deve ter seu envio de acordo com o calendário de férias publicado pelo RH no começo do ano.
1	$u_1$	Ah sim, obrigada.
1	$a_1$	Por nada.
2	$u_2$	Oi.
2	$a_2$	Olá, boa tarde.
2	$u_2$	Quero fazer uma pergunta.
2	$a_2$	Pois não?
2	$u_2$	Quando eu tenho direito a tirar férias?
2	$a_2$	Você já tirou férias depois que ocupou esse cargo?
2	$u_2$	Não.
2	$a_2$	Suas férias vencem após 12 meses da data do efetivo exercício.
3	$u_3$	Olá.
3	$a_1$	Olá.
3	$u_3$	Gostaria de fazer uma pergunta?
3	$u_3$	Eu já posso tirar licença?
3	$a_1$	Qual licença? Existem licenças como prêmio, saúde e especial.
3	$u_3$	Prêmio.
3	$a_1$	Você deve ter 5 períodos aquisitivos consecutivos.
3	$u_3$	Como assim 'por períodos aquisitivos'?
3	$a_1$	Após os 12 meses iniciais de efetivo exercício.
3	$u_3$	Quantos meses de licenças posso tirar?
3	$a_1$	Até 3 meses.

A segunda atividade que compõe o método é a de Aprender, que tem como objetivo construir automaticamente uma folksonomia a partir do corpus de diálogos e do pré-processamento realizado na atividade de Pré-processar. A atividade de Aprender possui as etapas de Obter Rótulos da Folksonomia, Obter Recursos da Folksonomia, Obter Relacionamentos entre Rótulos, Obter Atendentes da Folksonomia, Obter Usuários da Folksonomia, e Gerar Folksonomia, conforme mostra a Figura 15. A seguir serão descritas cada uma das etapas que compõem as atividades de Pré-processar e de Aprender.

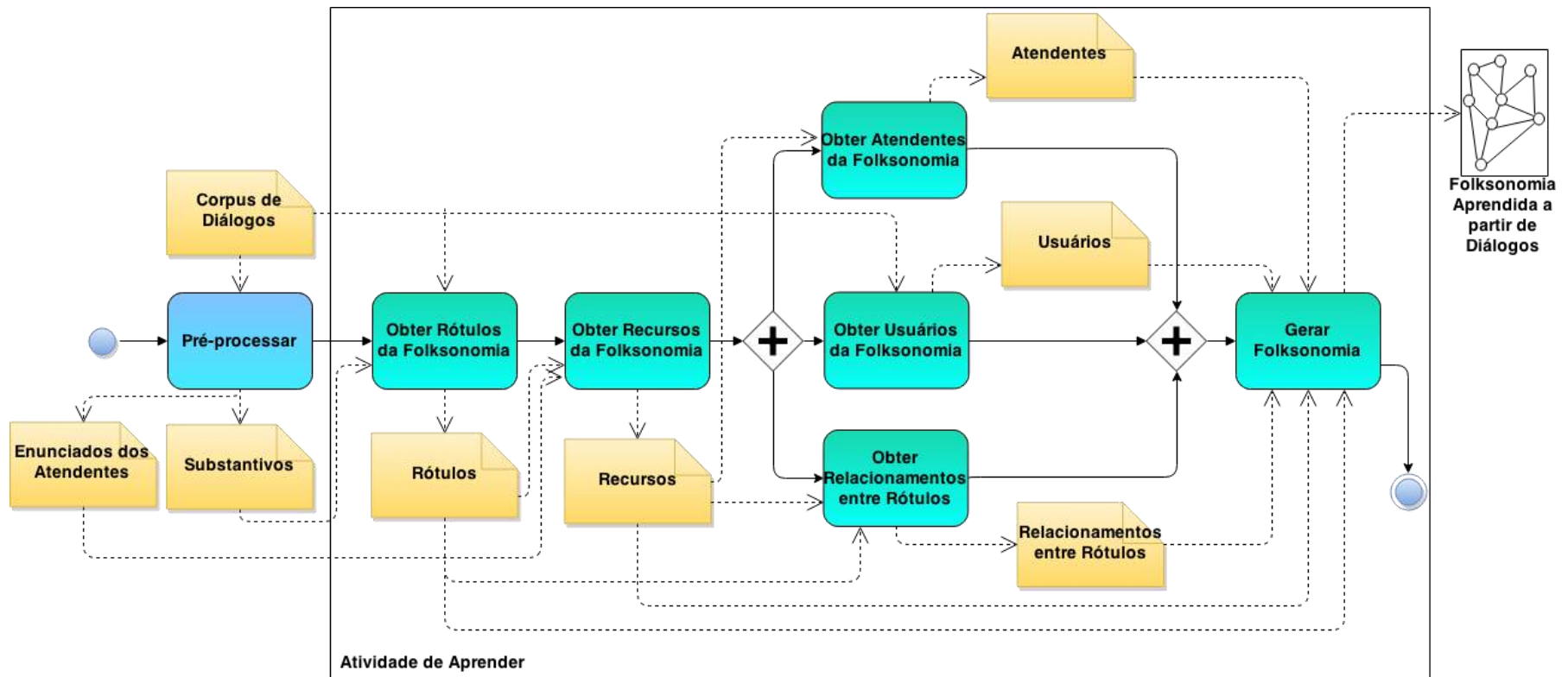


Figura 15 - Etapas da atividade de Aprender que compõe o método proposto.

### 5.2.1. Etapas da Atividade de Pré-processar

Nesta seção serão descritas as etapas de Extrair Enunciados dos Atendentes, Extrair Substantivos, e Remover Substantivos Duplicados, pertencentes à atividade de Pré-processar do método proposto.

**Etapa de Extrair Enunciados dos Atendentes:** esta etapa tem como objetivo receber o corpus de entrada e extrair dele apenas os enunciados dos atendentes, pois conforme apresentado na seção 5.1, existe a suposição de que os atendentes possuem conhecimento pleno do domínio dos diálogos. O principal intuito dessa “filtragem” em relação aos enunciados dos diálogos é encaminhar para as próximas etapas do método, somente os enunciados que representam o domínio em questão. Deste modo, a folksonomia aprendida representará o domínio do corpus de diálogos que foi selecionado. Formalmente, pode-se representar a obtenção do conjunto dos enunciados dos atendentes como um predicado unário  $Enu = \forall a, d, enu \text{ } Enu(enu) \rightarrow (a, d, enu) \in A \times D \times Enu$ , sendo  $enu$  um enunciado do atendente  $a$  num diálogo  $d$  pertencente ao corpus  $D$ .

O Quadro 4 mostra o resultado da etapa de extração dos enunciados dos atendentes, em relação aos diálogos do Quadro 3. Após a obtenção dos enunciados dos atendentes, a próxima etapa do método FolksDialogue é extrair os substantivos que compõem tais enunciados. A seguir é descrita a Etapa de Extrair Substantivos dos enunciados dos atendentes.

**Etapa de Extrair Substantivos:** Esta etapa da atividade de Pré-processar visa extrair os substantivos dos enunciados dos atendentes que foram obtidos na etapa anterior. O objetivo dessa extração é dar início ao processo de obtenção dos substantivos que posteriormente serão convertidos nos rótulos da folksonomia aprendida pelo método FolksDialogue. A identificação dos substantivos nos enunciados de  $Enu$  é realizada através de análise morfológica. Formalmente, os substantivos extraídos de  $Enu$  podem ser representados por um multiconjunto (o qual admite repetições em seus elementos)  $S$ :

$$S := \{\infty. sub : (sub \in enu) \wedge (enu \in Enu)\}$$

Sendo:

- $enu$  um dado enunciado de um atendente, pertencente ao conjunto  $Enu$ ;

- *sub* os substantivos pertencentes ao um dado enunciado *enu*;
- $\infty$  (infinito) a cardinalidade de *sub* em *S*.

**Quadro 4** - Resultado da Extração dos Enunciados dos Atendentes.

ID do Diálogo	Tipo de Interlocutor	Enunciado
1	$a_1$	Boa tarde.
1	$a_1$	Pois não, pode fazer.
1	$a_1$	Precisa enviar uma cópia do documento intitulado Planilha de Férias indicando se você deseja converter 10 dias de férias em pecúnia ou não e indicando o início do período de fruição com a sua assinatura e a assinatura da sua chefia.
1	$a_1$	Esse documento deve ter seu envio de acordo com o calendário de férias publicado pelo RH no começo do ano.
1		Por nada.
2	$a_2$	Olá, boa tarde.
2	$a_2$	Pois não?
2	$a_2$	Você já tirou férias depois que ocupou esse cargo?
2	$a_2$	Suas férias vencem após 12 meses da data do efetivo exercício.
3	$a_1$	Olá.
3	$a_1$	Qual licença? Existem licenças como prêmio, saúde e especial.
3	$a_1$	Você deve ter 5 períodos aquisitivos consecutivos.
3	$a_1$	Após os 12 meses iniciais de efetivo exercício.
3	$a_1$	Até 3 meses.

O Quadro 5 mostra em destaque os substantivos que foram extraídos dos enunciados dos diálogos pertencentes ao Quadro 4. Todos os substantivos extraídos e que estão em destaque no Quadro 5 podem ser vistos de modo específico no Quadro 6.

Com os substantivos extraídos dos enunciados dos atendentes, o próximo passo da atividade de Pré-processar é remover os que estão duplicados. O processo de remoção dos substantivos duplicados é apresentado na etapa a seguir.

**Etapa de Remover Substantivos Duplicados:** a partir do multiconjunto de substantivos *S* obtido com a execução da etapa anterior, o intuito desta etapa é remover os que estão duplicados. Para um substantivo ser considerado duplicado, é independente o fato de ele estar em maiúsculo ou minúsculo, e/ou no singular ou plural. Por exemplo, os substantivos *Mês* e *meses* são considerados iguais, portanto um deles deve ser removido, pois está duplicando o outro. A saída final desta etapa e da atividade de Pré-processar é um conjunto de substantivos (*Ls*) **únicos, em minúsculo e no singular**. Formalmente, *Ls* pode ser

representado pelo conjunto  $Ls := \{sub : (sub \in S)\}$ , sendo: sub um dado substantivo pertencente ao conjunto  $S$ .

**Quadro 5** - Resultado da Extração dos Substantivos dos Enunciados dos Atendentes.

ID do Diálogo	Tipo de Interlocutor	Enunciado
1	$a_1$	Boa tarde.
1	$a_1$	Pois não, pode fazer.
1	$a_1$	Precisa enviar uma cópia do documento intitulado Planilha de Férias indicando se você deseja converter 10 dias de férias em pecúnia ou não e indicando o início do período de fruição com a sua assinatura e a assinatura da sua chefia.
1	$a_1$	Esse documento deve ter seu envio de acordo com o calendário de férias publicado pelo RH no começo do ano.
1	$a_1$	Por nada.
2	$a_2$	Olá, boa tarde.
2	$a_2$	Pois não?
2	$a_2$	Você já tirou férias depois que ocupou esse cargo?
2	$a_2$	Suas férias vencem após 12 meses da data do efetivo exercício.
3	$a_1$	Olá.
3	$a_1$	Qual licença? Existem licenças como prêmio, saúde e especial.
3	$a_1$	Você deve ter 5 períodos aquisitivos consecutivos.
3	$a_1$	Após os 12 meses iniciais de efetivo exercício.
3	$a_1$	Até 3 meses.

**Quadro 6** - Substantivos extraídos dos enunciados dos atendentes.

Substantivos Extraídos
cópia – documento – Planilha – Férias – dias – férias - pecúnia - início - período - fruição - assinatura - assinatura – chefia - documento - acordo - calendário - férias – ano - férias – cargo - férias - meses – data – exercício - Licença - licenças – prêmio - saúde – especial – períodos - meses – exercício - meses

No Quadro 7 são mostrados os substantivos provenientes da etapa anterior (Quadro 6), que ficaram remanescentes após a execução da etapa de remoção de duplicados. O Quadro 7 representa  $Ls$  gerado como saída da atividade de Pré-processar.

**Quadro 7** - Conjunto de substantivos resultante da remoção de duplicados.

Conjunto de Substantivos
cópia – documento – planilha – férias – dia - pecúnia - início - período - fruição - assinatura – chefia - acordo - calendário – ano – cargo – mês – data – exercício - licença – prêmio - saúde – especial

A partir do conjunto de substantivos resultante da atividade de Pré-processar, é iniciada a atividade de Aprender do método proposto. A seguir são descritas as etapas da atividade de Aprender do método FolksDialogue.

### 5.2.2. Etapas da Atividade de Aprender

Nesta subseção são descritas as etapas de Obter Rótulos da Folksonomia, Obter Recursos da Folksonomia, Obter Relacionamentos entre os Rótulos, Obter Atendentes da Folksonomia, Obter Usuários da Folksonomia, e Gerar Folksonomia, pertencentes à atividade de Aprender do método proposto.

**Etapa de Obter Rótulos da Folksonomia:** esta etapa visa selecionar os substantivos do conjunto  $Ls$  gerado pela saída da atividade de Pré-processar para serem os Rótulos que comporão a folksonomia que está sendo aprendida. Para isso, realiza-se nesta etapa o que se denomina de *ranqueamento dos substantivos*. O objetivo do ranqueamento é obter o IDF<sup>7</sup> (*Inverse Document Frequency*) com que cada um dos substantivos de  $Ls$  aparece nos diálogos do corpus de diálogos. Os substantivos que tiverem seu IDF (chamado de IDFsub) abaixo de uma frequência de corte  $fc_1$  (tal que  $fc_1 \in \mathbb{N}^*$ ), calculada pela Equação (5.1), são descartados. Isto se deve ao fato de que para o contexto desta pesquisa, o IDF representa a importância de cada um dos substantivos de  $Ls$  no corpus de diálogos. Além disso, supõe-se que os substantivos com importância menor que a representada por  $fc_1$  podem não fazer parte do domínio em questão. Deste modo, se eles fossem incorporados na folksonomia estariam divergindo a representação de tal domínio. Os substantivos que não foram descartados após a aplicação de  $fc_1$  são considerados como pertencentes ao domínio, e são os rótulos da folksonomia que está sendo obtida.

$$fc_1 = \frac{\sum_{i=1}^{|Ls|} IDFsub_i}{|Ls|} \quad (5.1)$$

Formalmente a obtenção do conjunto  $T$  de rótulos da folksonomia, pode ser descrita como:  $T = \{sub : (sub \in Ls) \wedge (IDFsub \geq fc_1)\}$ .

O Quadro 8 mostra o ranqueamento dos substantivos de acordo com a aparição deles

<sup>7</sup> Maiores informações em: <<http://www.tfidf.com/>>

nos enunciados dos atendentes do Quadro 4. A fim de exemplificação, supõe-se que todos os substantivos do Quadro 7 possuem seus  $IDF_{sub}$  acima ou igual a  $fc_1$ , ou seja, todos pertencem ao conjunto  $T$ .

**Quadro 8 - Ranqueamento dos Substantivos.**

Substantivos	$IDF_{sub}$
cópia	0,32
documento	0,56
planilha	0,32
férias	0,86
dia	0,76
pecúnia	0,32
início	0,32
período	0,32
fruição	0,32
assinatura	0,56
chefia	0,32
acordo	0,32
calendário	0,32
ano	0,32
cargo	0,32
mês	0,75
data	0,56
exercício	0,68
licença	0,56
prêmio	0,32
saúde	0,32
especial	0,32

Os rótulos da folksonomia são a saída desta etapa da atividade de Aprender. Após a aquisição dos rótulos, ocorre a etapa de Obter Recursos da Folksonomia (apresentada a seguir).

**Etapa de Obter Recursos da Folksonomia:** esta etapa tem como intuito obter os recursos da folksonomia que está sendo aprendida. Para a seleção dos enunciados dos atendentes que posteriormente se tornarão os recursos da folksonomia, são utilizados os rótulos do conjunto  $T$  gerados na saída da etapa anterior. Dado ao fato de que esses rótulos são substantivos pertencentes ao domínio em questão, nesta etapa são verificados quantos dos substantivos (“*sub*”) de um dado enunciado pertencem ao grupo desses rótulos. O objetivo disso é verificar quais enunciados dos atendentes (“*enu*”) pertencem ao domínio, para que depois eles possam ser adotados como recursos. Para cada enunciado dos atendentes é calculada uma Taxa de Inclusão  $p_{enu}$  (definida de acordo com a Equação (5.2)). Tal taxa mede o percentual de substantivos de um dado enunciado que são rótulos da folksonomia, ou seja, que pertencem ao domínio. Os enunciados que tiverem suas Taxas de Inclusão  $p_{enu}$ , maiores ou iguais a uma Taxa de Inclusão  $p_1$  (tal que  $p_1 \in \mathbb{R}_+$ ), definida pela Equação (5.3), serão

adotados como recursos da folksonomia que está sendo obtida pelo método proposto. Vale destacar, que somente os enunciados que forem compostos por mais que 1 (um) substantivo são aplicados a Equação (5.2). Enunciados com apenas 1 (um) substantivo podem ser “genéricos”, e conseqüentemente não agregar nenhum conhecimento para um dado domínio. Um exemplo de enunciado “genérico” pode ser: “Bom dia, como vai você?”.

$$p_{enu} = \left( \frac{|\{sub : (sub \in enu) \wedge (sub \in T) \wedge (enu \in Enu)\}|}{|\{sub : (sub \in enu) \wedge (enu \in Enu)\}|} \right) \times 100 \quad (5.2)$$

Sendo:

- $enu \in Enu$ , um dado enunciado que um atendente gerou;
- $sub$  um substantivo de  $enu$ ;
- $T$  o conjunto finito de rótulos da folksonomia proposta pelo método.

$$p_1 = \frac{\sum_{i=1}^{|Enu|} p_{enu_i}}{|Enu|} \quad (5.3)$$

Deste modo, a obtenção do conjunto  $R$  dos recursos da folksonomia pode ser formalmente representada por:

$$R := \{enu : (enu \in Enu) \wedge (p_{enu} \geq p_1)\}, \text{ com } p_{enu}, p_1 \in \mathbb{R}_+$$

Sendo:

- $enu$  um dado enunciado de um atendente, pertencente a  $Enu$ ;
- $p_{enu}$  a Taxa de Inclusão de substantivos de  $enu$  que pertencem ao conjunto dos rótulos  $T$  da folksonomia;
- $p_1$  uma Taxa de Inclusão definida de forma empírica.

Com o intuito de exemplificar esta etapa do método e visando facilitar a compreensão, adota-se que todos os enunciados dos atendentes do Quadro 4, que possuem mais que 1 (um) substantivo, têm suas respectivas Taxas de Inclusão  $p_{enu}$  maiores ou iguais a  $p_1$ . No Quadro 9 é possível ver os enunciados do Quadro 4 que formam o conjunto  $R$  de recursos da folksonomia do método proposto.



**Quadro 9** - Enunciados dos Atendentes que são os Recursos da folksonomia Proposta.

ID do Diálogo	Tipo de Interlocutor	Enunciado
1	$a_1$	Precisa enviar uma cópia do documento intitulado Planilha de Férias indicando se você deseja converter 10 dias de férias em pecúnia ou não e indicando o início do período de fruição com a sua assinatura e a assinatura da sua chefia.
1	$a_1$	Esse documento deve ter seu envio de acordo com o calendário de férias publicado pelo RH no começo do ano.
2	$a_2$	Você já tirou férias depois que ocupou esse cargo?
2	$a_2$	Suas férias vencem após 12 meses da data do efetivo exercício.
3	$a_1$	Qual licença? Existem licenças como prêmio, saúde e especial.
3	$a_1$	Você deve ter 5 períodos aquisitivos consecutivos.
3	$a_1$	Após os 12 meses iniciais de efetivo exercício.

Finalizada a Etapa de Obter Recursos da folksonomia, três etapas ocorrem em paralelo: Obter Relacionamentos entre os Rótulos, Obter Atendentes da Folksonomias, e Obter Usuários da Folksonomia. A seguir é descrita a etapa de Obter Relacionamentos entre os Rótulos da folksonomia que está sendo obtida pelo método FolksDialogue.

**Etapa de Obter Relacionamentos entre os Rótulos:** o objetivo desta etapa é obter os possíveis relacionamentos (coocorrência) entre os rótulos que compõem a folksonomia que está sendo aprendida pelo método. Para isso, inicialmente geram-se todos os possíveis pares de rótulos a partir do conjunto  $T$ , isto é, combinação  $C_2^k$  sendo  $k$  o número de rótulos de  $T$ . Exemplificando, para  $k = 100$  o número de pares de rótulos a serem gerados será de 4.950. Na sequência, para cada par de rótulos gerado é obtida a frequência (“ $fpar$ ”) com que os dois termos que compõem o par, aparecem juntos rotulando os mesmos recursos do conjunto  $R$  da folksonomia. Ou seja, tal frequência é o número de recursos que eles rotularam juntos. Os pares de rótulos que possuem suas frequências  $fpar$  menores que uma frequência  $fc_2$  (tal que  $fc_2 \in \mathbb{N}^*$ ), calculada pela Equação (5.4), são descartados. Em oposição, os pares que possuem suas frequências maiores ou iguais a  $fc_2$ , terão um relacionamento  $b \in B$  (conjunto de relacionamento entre rótulos) entre os seus rótulos.

$$fc_2 = \frac{\sum_{i=1}^{|C_2^k|} fpar_i}{|C_2^k|} \quad (5.4)$$

Formalmente, pode-se definir o processo de obtenção do conjunto  $B$  de relacionamentos entre os rótulos da folksonomia como:  $B := \{b\}$ , sendo  $b$  um relacionamento entre dois dados

rótulos, e  $b$  pode ser dado por:  $\forall a, u, r, t_A, t_B b(t_A, t_B) \rightarrow ((a, t_A, r, u) \in Y' \wedge (a, t_B, r, u) \in Y' \wedge t_A \neq t_B \wedge fpar \geq fc_2)$ .

Vale ressaltar que cada um dos relacionamentos de  $B$  possui um determinado peso  $w$ , o qual, para esta pesquisa, é o número de diálogos em que os rótulos de um dado relacionamento  $b$  aparecem juntos, conforme já descrito e formalizado na Seção 5.1.

No Quadro 10 é possível ver os pares de rótulos formados a partir do conjunto  $T$ , as frequências de aparição de cada par nos recursos de  $R$ , e o peso  $w$  (número de diálogos em que os componentes de cada par aparecem juntos). Dados os 21 rótulos de  $T$ , foram gerados 210 pares distintos (combinação  $C_2^{21}$ ), porém com o intuito de facilitar a visualização, os pares de rótulos que tiveram suas frequências nulas foram omitidos do Quadro 10, resultando no final em 74 pares. A fim de exemplificar esta etapa do método, considera-se que todos os pares de rótulos do Quadro 10 possuem suas frequências iguais ou acima de uma determinada frequência de corte  $fc_2$ , assim, os rótulos que compõem cada um desses pares possuirão um relacionamento entre eles.

Na atividade de Aprender, em paralelo com a etapa de Obter Relacionamentos entre os Rótulos, ocorrem as etapas de Obter Atendentes da Folksonomia e de Obter Usuários da Folksonomia. A seguir será apresentada a Etapa de Obter Atendentes da Folksonomia.

**Etapa de Obter Atendentes da Folksonomia:** esta etapa visa obter o conjunto dos atendentes, os quais são os “usuários” da folksonomia proposta nesta pesquisa, conforme descrito na seção 5.1. Para cada recurso (enunciado de um atendente) do conjunto  $R$  da folksonomia, são extraídos todos os atendentes distintos “ $a$ ”, e com eles é formado o conjunto  $A$  da folksonomia do FolksDialogue. Formalmente, a obtenção do conjunto  $A$  de atendentes pode ser definida como:  $A := \{a : (a \in r) \wedge (r \in R)\}$ , sendo  $r$  um dado recurso do conjunto  $R$  de recursos da folksonomia;  $a$  o atendente que produziu  $r$ .

Exemplificando esta etapa, dado os enunciados dos atendentes do Quadro 4, o resultado da obtenção do conjunto de atendentes é  $A := \{a_1, a_2\}$ . Na sequência é descrita a etapa de Obter Usuários da Folksonomia, que é a última etapa que ocorre em paralelo com a de Obter Atendentes da Folksonomia.

Quadro 10 - Pares de rótulos: com suas frequências e seus pesos  $w$ .

Par de Rótulos		Frequência do Par nos Recursos	Peso $w$	Par de Rótulos		Frequência do Par nos Recursos	Peso $w$
cópia	documento	1	1	férias	acordo	1	1
cópia	planilha	1	1	férias	calendário	1	1
cópia	férias	1	1	férias	ano	1	1
cópia	dia	1	1	férias	data	1	1
cópia	pecúnia	1	1	férias	exercício	1	1
cópia	início	1	1	dia	pecúnia	1	1
cópia	período	1	1	dia	início	1	1
cópia	fruição	1	1	dia	período	1	1
cópia	assinatura	1	1	dia	fruição	1	1
cópia	chefia	1	1	dia	assinatura	1	1
documento	planilha	1	1	dia	chefia	1	1
documento	férias	2	1	pecúnia	início	1	1
documento	dia	1	1	pecúnia	período	1	1
documento	pecúnia	1	1	pecúnia	fruição	1	1
documento	início	1	1	pecúnia	assinatura	1	1
documento	período	1	1	pecúnia	chefia	1	1
documento	fruição	1	1	início	período	1	1
documento	assinatura	1	1	início	fruição	1	1
documento	chefia	1	1	início	assinatura	1	1
documento	acordo	1	1	início	chefia	1	1
documento	calendário	1	1	período	fruição	1	1
documento	ano	1	1	período	assinatura	1	1
planilha	férias	1	1	período	chefia	1	1
planilha	dia	1	1	fruição	assinatura	1	1
planilha	pecúnia	1	1	fruição	chefia	1	1
planilha	início	1	1	assinatura	chefia	1	1
planilha	período	1	1	acordo	calendário	1	1
planilha	fruição	1	1	acordo	ano	1	1
planilha	assinatura	1	1	calendário	ano	1	1
planilha	chefia	1	1	mes	data	1	1
férias	dia	1	1	mes	exercício	1	1
férias	pecúnia	1	1	data	exercício	1	1
férias	início	1	1	licença	prêmio	1	1
férias	período	1	1	licença	saúde	1	1
férias	fruição	1	1	licença	especial	1	1
férias	assinatura	1	1	prêmio	saúde	1	1
férias	chefia	1	1	prêmio	especial	1	1
férias	cargo	1	1	saúde	especial	1	1
férias	mes	1	1				

**Etapa de Obter Usuários da Folksonomia:** esta etapa descreve a obtenção do conjunto  $U$  de usuários. O conjunto  $U$  dos interlocutores do tipo usuário (que buscaram auxílio com os atendentes) é obtido extraindo-se a partir do corpus de diálogos de entrada, todos os interlocutores do tipo “ $u$ ”. Formalmente, a obtenção do conjunto  $U$  pode ser definida como:  $U := \{u : (u \in d) \wedge (d \in D)\}$ , sendo  $d$  um dado diálogo pertencente ao corpus  $D$ .

Exemplificando esta etapa do método FolksDialogue, dado o corpus de diálogos do Quadro 3, o resultado da obtenção do conjunto  $U$  dos usuários é  $U := \{u_1, u_2, u_3\}$ .

Por fim, com o conjunto  $U$  de usuários obtido, e com os conjuntos  $A$ ,  $T$ ,  $R$  e  $B$  de atendentes, rótulos, recursos, e relacionamento entre rótulos, respectivamente, adquiridos nas etapas anteriores, se dá início a Etapa de Gerar Folksonomia.

**Etapa de Gerar Folksonomia:** esta etapa tem como intuito gerar a estrutura final da folksonomia, proposta por esta pesquisa. Dados os conjuntos  $A$ ,  $T$ ,  $R$ ,  $U$  e  $B$ , obtidos nas etapas anteriores da atividade de Aprender, realiza-se a devida ligação entre os elementos pertencentes a esses conjuntos, com base na relação quaternária  $Y'$  (descrita na seção 5.1). Para cada elemento do conjunto  $A$  dos atendentes, os quais são os “usuários” da folksonomia proposta, é extraída sua personomia a partir de  $Y'$ . O conjunto das personomias dos atendentes representa a folksonomia do FolksDialogue. Ou seja,  $F := \{P_a\}$ , sendo:  $P_a := (T_a, R_a, I_a, I_a)$  with  $I_a := \forall t, r, u I_a(t, r, u) \rightarrow (a, t, r, u) \in Y'$ . Posteriormente, os rótulos das personomias são interligados através dos relacionamentos existentes no conjunto  $B$  de relacionamentos.

Para o exemplo que está sendo desenvolvido ao longo do capítulo, o resultado da folksonomia obtida pelo método é apresentado na Figura 16. Com o intuito de melhorar a visualização, os relacionamentos entre os rótulos foram omitidos, porém na Figura 17, que mostra em detalhes a personomia do atendente  $a_2$ , é possível ver os relacionamentos existentes entre os rótulos, juntamente com seus respectivos pesos. Em ambas as figuras, os recursos estão representados como  $r_1, r_2, r_3, r_4, r_5, r_6, r_7, r_8$ , a razão disso, foi novamente para auxiliar na visualização. Os recursos estão representados da seguinte forma:

- $r_1$  corresponde ao recurso: “Precisa enviar uma cópia do documento intitulado Planilha de Férias indicando se você deseja converter 10 dias de férias em pecúnia ou não e indicando o início do período de fruição com a sua assinatura e a assinatura da sua chefia.”
- $r_2$  corresponde ao recurso: “Esse documento deve ter seu envio de acordo com o

calendário de férias publicado pelo RH no começo do ano.”

- $r_3$  corresponde ao recurso: “Você já tirou férias depois que ocupou esse cargo?”
- $r_4$  corresponde ao recurso: “Suas férias vencem após 12 meses da data do efetivo exercício.”
- $r_5$  corresponde ao recurso: “Qual licença? Existem licenças como prêmio, saúde e especial.”
- $r_6$  corresponde ao recurso: “Você deve ter 5 períodos aquisitivos consecutivos.”
- $r_7$  corresponde ao recurso: “Após os 12 meses iniciais de efetivo exercício.”
- $r_8$  corresponde ao recurso: “Até 3 meses.”

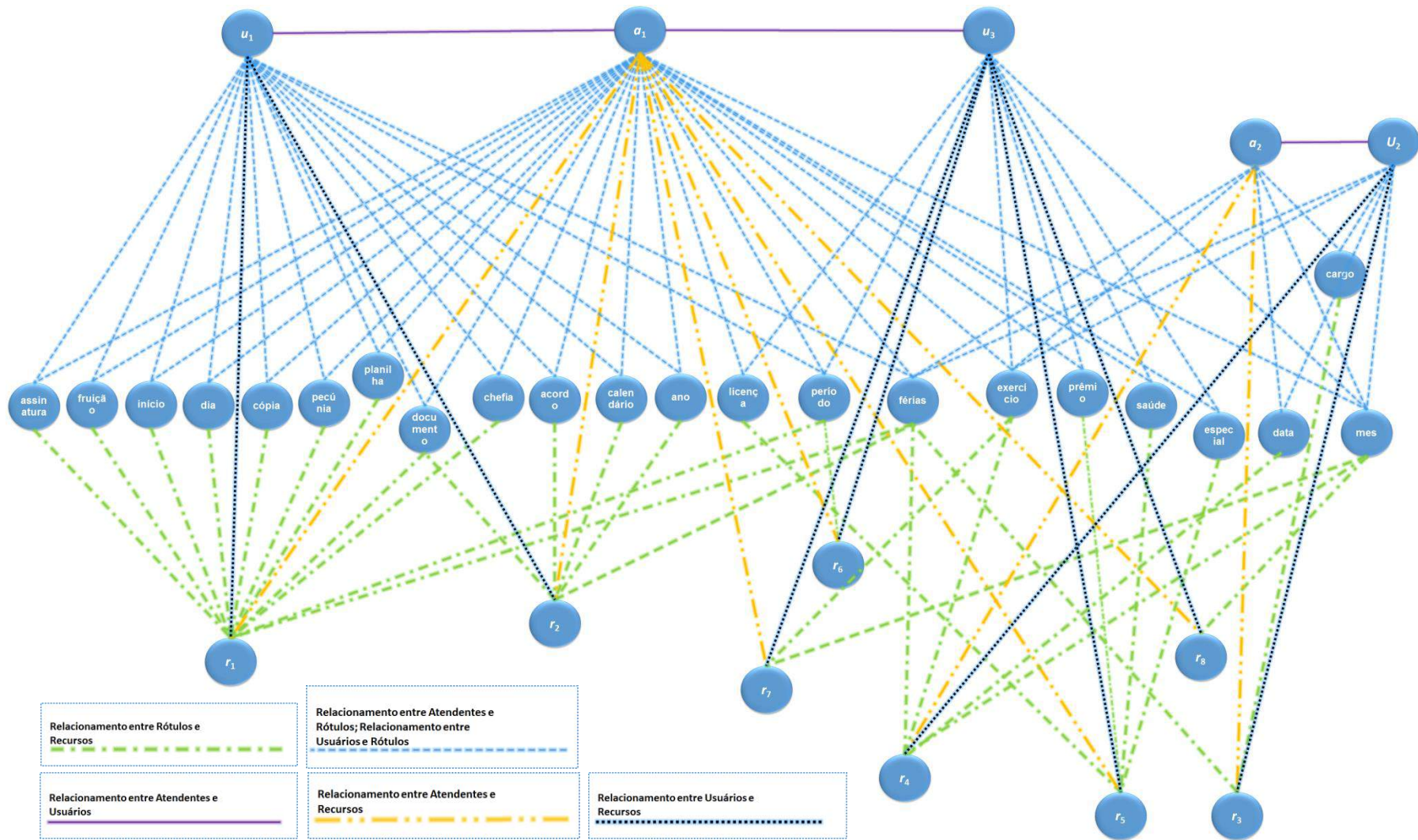


Figura 16 - Folksonomia aprendida pelo método FolksDialogue.

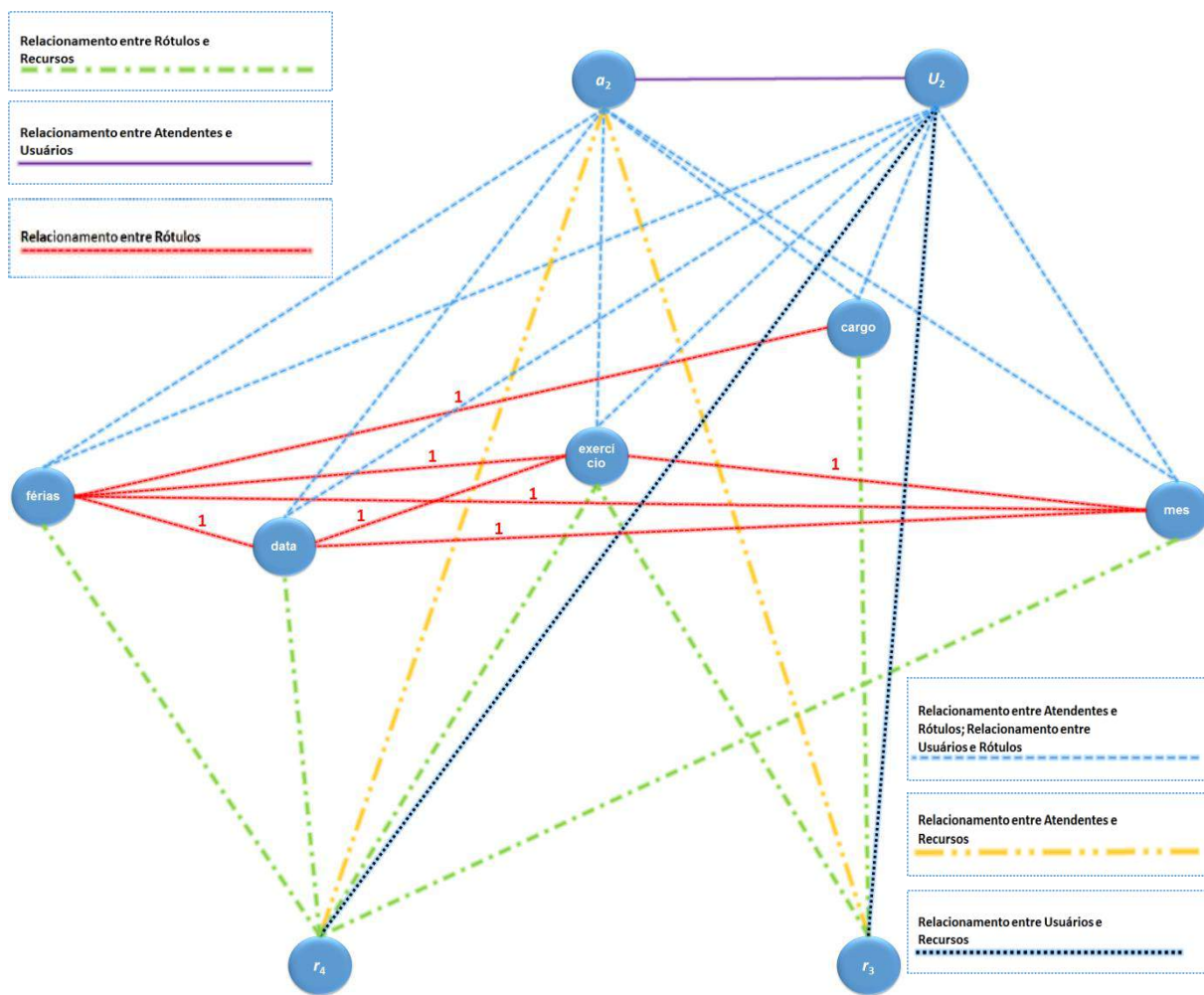


Figura 17 - Personomia do atendente  $a_2$ , extraída da folksonomia aprendida pelo método FolksDialogue.

### 5.3. Conclusão

Como visto neste capítulo, o método para o aprendizado de folksonomias a partir de diálogos engloba uma série de atividades e etapas que devem ser executadas e aplicadas ordenadamente para obtenção das folksonomias. Assim, no próximo capítulo são apresentadas as formas com que são implementadas as atividades e as etapas que compõem o método proposto.

# Capítulo 6

## Procedimentos Metodológicos

Neste capítulo é descrito o processo de implementação das atividades e de suas respectivas etapas que compõem o método proposto. Também é apresentada a forma de avaliação do método e como é realizada a detecção de tendências nas folksonomias aprendidas.

O processo de implementação ocorre através de um protótipo computacional, que tem como objetivo validar o método FolksDialogue. O protótipo computacional é representado pelo processo de implementação das atividades de Pré-processar e de Aprender do método.

Este capítulo está organizado da seguinte forma: inicialmente na seção 6.1 é apresentado o corpus de diálogos utilizado como entrada do método, e na seção 6.2 é descrito o protótipo computacional que implementa as atividades de Pré-processar e de Aprender do método. Em seguida na seção 6.3 é descrita a avaliação do método proposto, e após isso na seção 6.4 é apresentado um experimento para a detecção de tendências nas folksonomias geradas pelo método FolksDialogue.

### 6.1. Corpus de Diálogos

O corpus de diálogos disponível utilizado como entrada do método FolksDialogue foi obtido junto a uma prefeitura do estado do Paraná. O corpus é composto de 901 diálogos textuais em português do Brasil, os quais foram produzidos durante os anos de 2006 a 2009. Além disso, os diálogos do corpus estão dispostos em ordem cronológica.

O conteúdo dos diálogos (o domínio) é referente ao departamento de recursos



humanos, no qual interlocutores humanos dialogam sobre questões como: aposentadoria, direitos de ordem geral, estágio probatório, férias, licenças e provimento.

Cada diálogo do corpus possui as seguintes características:

- é composto por um identificador único;
- é formado por mensagens (enunciados) resultantes do *chat* entre interlocutores do tipo atendente e do tipo usuário;
- os enunciados estão identificados de acordo com o tipo de interlocutor que os gerou.

Computacionalmente, o corpus de diálogos está armazenado numa tabela  $X_I$  do banco de dados *PostgreSQL*<sup>8</sup>. No Quadro 11 é possível ver exemplos de diálogos pertencentes ao corpus e o formato em que eles estão dispostos no banco de dados. Nesse quadro a coluna **ID do Diálogo** é um número de identificação único que representa cada um dos diálogos do corpus. Os enunciados dos diálogos estão presentes no campo **Enunciado**. A identificação do interlocutor que gerou um dado enunciado é atribuída ao campo **Tipo de Interlocutor**, sendo: “ $a_i$ ” (tal que  $i \in \mathbb{N}^*$ ) a identificação para interlocutores do tipo atendente e “ $u_j$ ” (tal que  $j \in \mathbb{N}^*$ ) a identificação para interlocutores do tipo usuário.

Na próxima seção é apresentado o protótipo computacional, que implementa o método proposto nesta pesquisa.

## 6.2. Implementação do Protótipo Computacional

Esta seção descreve a implementação das atividades de Pré-processar e de Aprender do método *FolksDialogue*, através da linguagem Java.

**Quadro 11** - Exemplos de Diálogos do Corpus Utilizado nesta Pesquisa.

(continua)

ID do Diálogo	Tipo de Interlocutor	Enunciado
1	$u_1$	Oi, como vai?
1	$a_1$	Tudo bem, e você?
1	$u_1$	Tudo bem, gostaria de fazer uma pergunta.
1	$a_1$	Pois não, pode perguntar.

<sup>8</sup> Maiores informações em: <http://www.postgresql.org/>

Quadro 11 - Exemplos de Diálogos do Corpus Utilizado nesta Pesquisa.

(conclusão)

1	$u_1$	Existe cota p/ pessoas com deficiência na prefeitura?
1	$a_1$	Sim, existe. Ela é de cinco por cento dos cargos públicos da Administração Direta e Indireta.
1	$u_1$	valeu
1	$a_1$	Estamos à disposição.
2	$u_2$	Olá, tudo bem?
2	$a_2$	Tudo bem, posso ajudar?
2	$u_2$	Por que existem duas idades para a pessoa se aposentar por tempo de serviço, aos 25 e aos 30 anos, para mulher?
2	$a_2$	O tempo de 25 anos é um tempo mínimo que só dá direito à aposentadoria, mas não a receber proventos integrais. O tempo de 30 anos já é a idade certa, que dá direito a receber proventos integrais.
2	$u_2$	proventos?
2	$a_2$	é como o salário.
2	$u_2$	a sim, obrigada
2	$a_2$	Por nada.
3	$u_3$	Bom dia. Com quanto dos anos de trabalho posso pedir minha aposentadoria?
3	$a_1$	bom dia, conforme o Estatuto, com 35 anos de contribuição se homem e 30 de contribuição se mulher.
3	$u_3$	Obrigada.
3	$a_1$	Estamos à disposição.

### 6.2.1. Implementação da Atividade de Pré-processar

Para facilitar a compreensão, a seguir são detalhadas separadamente as formas de implementação de cada uma das etapas que compõem a atividade de Pré-processar. Inicialmente é descrita a implementação da Etapa de Extrair Enunciados dos Atendentes, em seguida a da Etapa de Extrair Substantivos, e por fim, da Etapa de Remover Substantivos Duplicados.

**Etapa de Extrair Enunciados dos Atendentes:** esta etapa visa extrair da tabela  $X_1$  do banco de dados, que possui o corpus de diálogos, somente os enunciados produzidos pelos atendentes. Para isso, criar-se uma função que realiza um comando *select* no banco de dados, extraindo da tabela  $X_1$  somente os enunciados que são do tipo “ $a_i$ ” (tal que  $i \in \mathbb{N}^*$ ) na coluna Tipo de Interlocutor. Os enunciados extraídos são armazenados juntamente com seus interlocutores e IDs dos diálogos, numa tabela  $X_2$  do banco de dados.

**Etapa de Extrair Substantivos:** o intuito desta etapa é extrair os substantivos dos enunciados da tabela  $X_2$  do banco de dados. Para a extração dos substantivos é utilizado o *parser* para a língua portuguesa do *CoGrOO*<sup>9</sup> (Corretor Gramatical Acoplável ao *OpenOffice.org*<sup>10</sup>) em sua versão 4.0.0. A escolha por esse *parser* se deu devido a ele possuir uma ampla biblioteca, a qual inclui um lematizador (descrito na próxima etapa), e ser de fácil utilização e integração com a linguagem Java. Após a extração dos substantivos eles são armazenados num vetor  $V_I$ .

**Etapa de Remover Substantivos Duplicados:** após a extração dos substantivos realizada na etapa anterior, o próximo passo é remover os que estão duplicados no vetor  $V_I$ . Conforme descrito no Capítulo 5, para um dado substantivo ser considerado duplicado independe o fato de ele estar em maiúsculo ou minúsculo, e no singular ou plural. O processo de remoção dos substantivos que estão duplicados em  $V_I$  é realizado com o auxílio do lematizador do *CoGrOO*, e com o uso da classe *Collator*<sup>11</sup> do Java. O lematizador tem como objetivo extrair apenas o lema dos substantivos de  $V_I$ , deste modo é possível ignorar o fato de eles estarem no singular ou no plural. Após a lematização dos substantivos é utilizada a classe *Collator*, que possui como uma de suas funções a comparação de *strings* de acordo com regras e caracteres específicos da língua em que elas foram escritas. Com isso, dois substantivos quaisquer, no âmbito de estarem em minúsculo ou maiúsculo, e com ou sem acento, passam a ser considerados iguais. Os substantivos considerados duplicados são removidos de  $V_I$ , e os que permanecerem remanescentes são armazenados num vetor  $L_S$ , que representa o conjunto de substantivos descrito no Capítulo 5, o qual é a saída da atividade de Pré-processar.

### 6.2.2. Implementação da Atividade de Aprender

A implementação das etapas da atividade de Aprender é apresentada na seguinte ordem: primeiramente se descreve a implementação da Etapa de Obter Rótulos da folksonomia, e em seguida a da Etapa de Obter Recursos da folksonomia é apresentada. Na

---

<sup>9</sup> <<http://cogroo.sourceforge.net/>>

<sup>10</sup> <<http://www.openoffice.org/>>

<sup>11</sup> <<http://docs.oracle.com/javase/7/docs/api/java/text/Collator.html/>>

sequência é descrita a da Etapa de Obter Relacionamentos entre os Rótulos, e posteriormente apresenta-se a implementação da Etapa de Obter Atendentes da folksonomia. Logo após, as implementações das Etapas de Obter Usuários da folksonomia e de Gerar folksonomia são apresentadas, respectivamente.

**Etapa de Obter Rótulos da Folksonomia:** esta etapa visa obter os rótulos da folksonomia que está sendo aprendida pelo método proposto. A partir do vetor  $Ls$  gerado na saída da atividade de Pré-processar é realizado o *ranqueamento dos substantivos* descrito no Capítulo 5. O objetivo do ranqueamento é obter para cada substantivo de  $Ls$ , a sua frequência de aparição no vetor  $V_1$ , o qual é composto por todos os substantivos dos enunciados dos atendentes da tabela  $X_2$ . Após a extração das frequências, cada substantivo de  $Ls$  é armazenado juntamente com a sua frequência num vetor  $V_3$ . Os substantivos de  $V_3$  que possuírem suas frequências maiores ou iguais a uma frequência de corte  $fc_1$  (Equação (5.1)) são considerados nesta pesquisa como pertencentes ao domínio do corpus de diálogos de entrada e são armazenados num vetor  $T$ . O vetor  $T$  representa o conjunto de rótulos da folksonomia proposta pelo método.

**Etapa de Obter Recursos da Folksonomia:** esta etapa tem como intuito selecionar os enunciados dos atendentes da tabela  $X_2$ , para serem os recursos da folksonomia obtida pelo FolksDialogue. A seleção de um dado enunciado ocorre verificando se o mesmo pertence ou não ao domínio do corpus de entrada. Essa verificação é realizada através do cálculo da Taxa de Inclusão estabelecida pela Equação (5.2), mostrada no Capítulo 5. O objetivo da Taxa de Inclusão é medir o percentual de substantivos de um dado enunciado que são rótulos da folksonomia. Para isso, inicialmente é codificada a Equação (5.2) como sendo uma função  $Penu$ , que recebe como entrada um dado enunciado da tabela  $X_2$ , juntamente com o vetor  $T$  de rótulos da folksonomia, e retorna como saída o valor calculado da Taxa de Inclusão. Dentro da função  $Penu$  é realizada uma verificação que testa se o número mínimo de substantivos que compõem um dado enunciado, recebido como parâmetro de entrada, é maior que 1 (um). Esta verificação ocorre com o auxílio do *parser* do CoGrOO, e tem como objetivo evitar os enunciados “genéricos” citados no Capítulo 5.

Após codificar a função  $Penu$ , todos os enunciados de  $X_2$  são submetidos a ela gerando como saída a Taxa de Inclusão de cada um. Na sequência, os enunciados de  $X_2$  são

armazenados num vetor  $R$  juntamente com seus respectivos “ID de Diálogo”, atendentes (provenientes da coluna Tipo de Interlocutor) e com suas Taxas de Inclusão calculadas. Os elementos de  $R$  que possuírem suas Taxas de Inclusão menores que uma Taxa de Inclusão  $p_1$  (Equação (5.3)), são removidos desse vetor. Os enunciados remanescentes em  $R$  são considerados, no âmbito desta pesquisa, como pertencentes ao domínio do corpus de entrada, e conseqüentemente são os recursos da folksonomia do método proposto.

**Etapa de Obter Relacionamentos entre Rótulos:** para a obtenção dos possíveis relacionamentos entre os rótulos da folksonomia é criada uma função que gera todos os possíveis pares de rótulos do vetor  $T$ . Em seguida, para cada par de rótulo gerado é obtida a frequência com que os termos que compõem o par, aparecem juntos nos recursos de  $R$ , ou seja, o número de recursos que eles rotulam juntos. Para a obtenção dessa frequência, inicialmente são extraídos os substantivos dos recursos de  $R$ . Após isso, realiza-se a comparação de igualdade entre os substantivos que compõem os pares e os substantivos extraídos dos recursos. Por fim, os pares de rótulos são armazenados juntamente com suas respectivas frequências num vetor  $B$ . Os pares que possuírem suas frequências menores que uma frequência  $fc_2$  (Equação (5.4)), são removidos de  $B$ . Os pares remanescentes nesse vetor irão possuir um relacionamento entre os seus rótulos.

O passo final desta etapa é atribuir os pesos nos relacionamentos dos pares de rótulos do vetor  $B$ . Conforme descrito no Capítulo 4 deste trabalho, os pesos dos possíveis relacionamentos entre os rótulos da folksonomia são para o âmbito desta pesquisa, o número de diálogos em que os rótulos de um dado par aparecem juntos. Deste modo, o peso do relacionamento para um dado par de  $B$  será o número de diálogos que possuem enunciados compostos pelos rótulos deste dado par. Neste caso, vale ressaltar conforme descrito na seção 5.1, que dois rótulos não precisam estar num mesmo enunciado de um dado diálogo, eles podem pertencer a enunciados distintos, porém devem pertencer ao mesmo diálogo. Além disso, os enunciados que são verificados se possuem os rótulos de  $B$ , são os recursos do vetor  $R$ . Já a identificação do diálogo a que um dado enunciado pertence é realizada pelo ID do Diálogo, que está associado dentro do próprio vetor  $R$  (conforme descrito na implementação da Etapa de Obter Recursos). O peso de cada relacionamento é armazenado juntamente com o mesmo no vetor  $B$ .

**Etapa de Obter Atendentes da Folksonomia:** o intuito desta etapa é obter os atendentes que farão parte da folksonomia proposta no âmbito desta pesquisa. Para isso, cria-se uma função que extrai todos os atendentes que estão associados aos seus recursos dentro do vetor  $R$ . Os atendentes extraídos são armazenados num vetor  $A$ .

**Etapa de Obter Usuários da Folksonomia:** esta etapa visa obter os usuários que vão compor a folksonomia do método FolksDialogue. Inicialmente, é criada uma função que extrai do corpus de entrada da tabela  $X_I$ , todos os interlocutores que são do tipo usuário (“ $u_j$ ”, tal que  $j \in \mathbb{N}^*$ ) na coluna Tipo do Interlocutor. Após a extração dos usuários, eles são armazenados num vetor  $U$ . Dentro deste mesmo vetor cada usuário é armazenado juntamente com o atendente com quem ele dialogou no corpus de entrada.

**Etapa de Gerar Folksonomia:** Com base nos vetores  $A$ ,  $T$ ,  $R$ ,  $U$  e  $B$ , de atendentes, rótulos, recursos, usuários, e relacionamentos entre rótulos, respectivamente, esta etapa tem como objetivo criar as devidas ligações entre os elementos desses vetores, a fim de se gerar a estrutura final da folksonomia do método proposto. Nesta pesquisa, a folksonomia obtida é representada computacionalmente através de um grafo quadripartido  $G_F := (V, Z, w_i)$ , sendo que:

- $V$  é o conjunto de vértices do grafo, formado pela união dos conjuntos  $A$ ,  $T$ ,  $R$ ,  $U$  da tupla  $\mathbb{F}$  da folksonomia proposta neste trabalho, ou seja,  $V = A \cup T \cup R \cup U$ ;
- $Z$  o conjunto de arestas, composto pelas atribuições da relação  $Y'$  da tupla  $\mathbb{F}$  e pelos relacionamentos que foram gerados entre os rótulos (conjunto  $B$ ), ou seja,  $Z$  é formado pelas relações abaixo:
  - $\forall a, t Z_R(a, t) \wedge Z_U(a, t) \rightarrow \exists r Y'(a, t, r, u) \wedge \exists u Y'(a, t, r, u)$  (existe um recurso e um usuário conectando um atendente a um rótulo)
  - $\forall a, r Z_T(a, r) \wedge Z_U(a, r) \rightarrow \exists t Y'(a, t, r, u) \wedge \exists u Y'(a, t, r, u)$  (existe um rótulo e um usuário conectando um atendente a um recurso)
  - $\forall a, u Z_T(a, u) \wedge Z_R(a, u) \rightarrow \exists t Y'(a, t, r, u) \wedge \exists r Y'(a, t, r, u)$  (existe um rótulo e um recurso conectando um atendente a um usuário)
  - $\forall u, t Z_A(u, t) \wedge Z_R(u, t) \rightarrow \exists a Y'(a, t, r, u) \wedge \exists r Y'(a, t, r, u)$  (existe um atendente e um recurso conectando um usuário a um rótulo)
  - $\forall u, r Z_A(u, r) \wedge Z_T(u, r) \rightarrow \exists a Y'(a, t, r, u) \wedge \exists t Y'(a, t, r, u)$  (um

atendente e um rótulo conectando um usuário a um recurso)

- $\forall a, u, r, t_A, t_B \mid (a, t_A, r, u) \in Y' \wedge (a, t_B, r, u) \in Y' \rightarrow b(t_A, t_B) \wedge (t_A \neq t_B)$   
(existe um relacionamento  $b$  entre dois rótulos  $t_A$  e  $t_B$ )

Assim,  $Z := Z_A \cup Z_T \cup Z_R \cup Z_U$ .

- $w_i$  os pesos dos relacionamentos do conjunto  $B$  da folksonomia.

Para se representar a folksonomia do método FolksDialogue através de um grafo é utilizado o banco de grafos Neo4j<sup>12</sup>. A escolha por esse banco foi devida ao fato dele ser de código aberto e possuir uma biblioteca para manipulação em Java, facilitando deste modo a integração com o restante do protótipo.

Conforme descrito no Capítulo 4, para a ligação entre os elementos dos conjuntos (no âmbito deste Capítulo, vetores)  $A, T, R, U$  é utilizada a relação quaternária  $Y'$ . Do ponto de vista da implementação do método proposto, a relação  $Y'$  está intrínseca nos elementos que compõem os vetores  $A, T, R, U$ . Assim, a implementação do processo de obtenção das personomias dos atendentes do vetor  $A$  pode ser representado pelo pseudocódigo do Quadro 12.

**Quadro 12** - Algoritmo para criação das personomias dos Atendentes.

(continua)

```
VARIÁVEIS
A: Vetor de Atendentes;
T: Vetor de Rótulos;
R: Vetor de Recursos;
U: Vetor de Usuários;
B: Vetor de Relacionamentos entre Rótulos;
SUB: Vetor de Substantivos;

INeo4j: Instância do Neo4j; {Instância do banco de grafos Neo4j}
IParser: Instância do Parser do CoGrOO; {Instância do Parser do CoGrOO}
IndVetA: Inteiro; {Índice para o Vetor A}
IndVetT: Inteiro; {Índice para o Vetor T}
IndVetR: Inteiro; {Índice para o Vetor R}
IndVetU: Inteiro; {Índice para o Vetor U}
IndVetB: Inteiro; {Índice para o Vetor B}
IndVetSUB: Inteiro; {Índice para o Vetor SUB}
FlagSUB: Inteiro; {Flag para detecção de substantivos}

INÍCIO

{Cria nós para os Rótulos}
PARA IndVetT ← 1 ATÉ Contador(T) FAÇA
  INeo4j.CriarNó(T[IndVetT]);
FIM PARA

{Cria nós para os Recursos}
PARA IndVetR ← 1 ATÉ Contador(R) FAÇA
  INeo4j.CriarNó(R[IndVetR].Recurso);
```

<sup>12</sup> Maiores informações em: <<http://www.neo4j.org/>>

Quadro 12 - Algoritmo para criação das personomias dos Atendentes.

(Continuação)

```

FIM PARA

{Criar Arestas entre Rótulos e Recursos}
PARA IndVetT ← 1 ATÉ Contador(T) FAÇA {Para cada Rótulo de T}

  PARA IndVetR ← 1 ATÉ Contador(R) FAÇA {Para cada Recurso de R}

    SUB ← IParser.ExtrairSubstantivos(R[IndVetR].Recurso); {Extrai os substantivos do
                                                                Recurso em questão}
    FlagSUB ← 0; {Zera o flag de detecção de substantivos}

    PARA IndVetSUB ← 1 ATÉ Contador(SUB) FAÇA {Para todos os substantivos do Recurso
                                                                em questão}

      SE T[IndVetT] = SUB[IndVetSUB] ENTAO
        FlagSUB ← 1; {encontrou um substantivo igual ao Rótulo em questão}
      FIM SE
    FIM PARA

    SE FlagSUB = 1 ENTAO
      {Cria aresta entre Rótulo e Recurso}
      INeo4j.CriarAresta(INeo4j.Nó(T[IndVetT]), INeo4j.Nó(R[IndVetR]));
    FIM SE
  FIM PARA
FIM PARA

{Cria a Personomia de cada Atendente}
PARA IndVetA ← 1 ATÉ Contador(A) FAÇA

  INeo4j.CriarNó(A[IndVetA]); {Cria nó para o Atendente em questão}

  {Ligar Atendente com seu subconjunto de Usuários}
  PARA IndVetU ← 1 ATÉ Contador(U) FAÇA {busca o subconjunto de Usuários do Atendente em
                                                                questão}

    SE U[IndVetU].Atendente = A[IndVetA] ENTAO {Se encontrou um Usuário do Atendente}
      INeo4j.CriarNó(U[IndVetU]); {Cria um nó para o Usuário}
      {Liga o nó do Atendente com o nó do Usuário}
      INeo4j.CriarAresta(INeo4j.Nó(A[IndVetA]), INeo4j.Nó(U[IndVetU].Usuário));
    FIM SE
  FIM PARA

  {Ligar Atendente com os Recursos produzidos por ele}
  PARA IndVetR ← 1 ATÉ Contador(R) FAÇA

    SE A[IndVetA] = R[IndVetR].Atendente ENTAO

      {Liga o nó do Atendente com o nó do Recurso}
      INeo4j.CriarAresta(INeo4j.Nó(A[IndVetA]), INeo4j.Nó(R[IndVetR].Recurso));
    FIM SE
  FIM PARA

  {Ligar Atendente com Rótulos}
  PARA IndVetT ← 1 ATÉ Contador(T) FAÇA {Para cada Rótulo de T}
    PARA IndVetR ← 1 ATÉ Contador(R) FAÇA {Para cada Recurso de R}
      {Se existe uma aresta entre o Rótulo e o Recurso em questão}
      SE INeo4j.ExisteAresta(INeo4j.Nó(T[IndVetT]), INeo4j.Nó(R[IndVetR])) ENTAO
        {Se existe uma aresta entre o Recurso e o Atendente em questão}
        SE INeo4j.ExisteAresta(INeo4j.Nó(R[IndVetR]), INeo4j.Nó(A[IndVetA])) ENTAO

          {Liga o nó do Atendente com o nó do Rótulo}
          INeo4j.CriarAresta(INeo4j.Nó(A[IndVetA]), INeo4j.Nó(T[IndVetT]));
        FIM SE
      FIM SE
    FIM PARA
  FIM PARA

  {Ligar Usuários do Atendente com os Rótulos}
  PARA IndVetU ← 1 ATÉ Contador(U) FAÇA {busca o subconjunto de Usuários do Atendente em
                                                                questão}

    SE U[IndVetU].Atendente = A[IndVetA] ENTAO {Se encontrou um usuário do Atendente}
      PARA IndVetR ← 1 ATÉ Contador(R) FAÇA {Para cada Recurso de R}

```





**Quadro 13** - Algoritmo para Ligação dos Rótulos das personomias dos Atendentes.

```

VARIÁVEIS
B: Vetor de Relacionamentos entre Rótulos;

INeo4j: Instância do Neo4j; {Instância do banco de grafos Neo4j}
IndVetB: Inteiro; {Índice para o Vetor B}

INÍCIO

  {Cria Aresta entre os Rótulos que possuem Relacionamento}
  PARA IndVetB ← 1 ATÉ Contador(B) FAÇA

      INeo4j.CriarAresta(INeo4j.Nó(B[IndVetB].RótuloA), INeo4j.Nó(B[IndVetB].RótuloB),
                          Peso(B[IndVetB].Peso));

  FIM PARA

FIM.

```

### 6.3. Avaliação

A avaliação do método proposto ocorre através de duas abordagens, uma denominada de **avaliação de característica** e a outra de **avaliação de teste de domínio**. A avaliação de característica tem como objetivo verificar se a estrutura da folksonomia do método FolksDialogue possui a característica de fenômeno de mundo-pequeno, descrita no Capítulo 2. O fato da estrutura da folksonomia obtida possuir tal característica pode ser uma das comprovações de que ela é genuinamente uma folksonomia. Por outro lado, a avaliação de teste de domínio visa medir a acurácia da folksonomia obtida por esta pesquisa, no processo de interpretação de enunciados de diálogos, verificando se eles pertencem ou não ao domínio representado que ela representa.

Nas subseções a seguir são detalhadas as abordagens de avaliação de característica e de teste de domínio do método proposto.

#### 6.3.1. Avaliação de Característica

Na busca por medidas de avaliação não foram encontradas métricas consensuais. Porém, foi encontrado um trabalho (CATTUTO et al., 2007), no qual os autores avaliaram a folksonomia através do aspecto de mundo-pequeno. O objetivo da avaliação de característica é verificar se a estrutura da folksonomia obtida pelo método possui a característica de mundo-

pequeno. Ou seja, é verificado se os nós da folksonomia estão fortemente agrupados (coeficiente de agrupamento elevado em relação a um grafo aleatório), e se o comprimento de caminho característico do grafo dela é pequeno (próximo ao de um grafo aleatório).

No âmbito da folksonomia gerada pelo método, o cálculo do comprimento de caminho característico pode ser obtido pela Equação (6.1), a qual é uma adaptação da Equação (2.1) do Capítulo 2. Além disso, seguindo a abordagem de (CATTUTO et al., 2007), neste trabalho na Equação (6.1) somente são calculadas as distâncias entre dois vértices que possuem um caminho entre si. Nesta pesquisa, o caminho mais curto entre dois nós do grafo da folksonomia obtida, se dá pelo menor número de arestas que conectam eles. Para o cálculo do caminho mais curto é utilizado o algoritmo *Breadth-First Search* (SEDGWICK e WAYNE, 2011), o motivo para a escolha desse algoritmo é o fato de ele já foi utilizado com sucesso para esse mesmo fim em (CATTUTO et al., 2007).

$$Q_F = \text{mediana}\{\text{média}\{d(v, x) : (v \in V) \wedge (x \in V - \{v\})\}\}, \quad (6.1)$$

sendo:  $V = A \cup T \cup R \cup U$ .

Para o cálculo do coeficiente de agrupamento de um dado nó da folksonomia obtida, realiza-se uma extensão da Equação (2.2) do coeficiente de agrupamento “*cliquishness*” apresentado por (CATTUTO et al., 2007), o qual foi proposto para o modelo tripartido de folksonomias. Deste modo, seja um recurso  $r$  (tal que  $r \in R$ ) um nó da folksonomia proposta através do modelo quadripartido desta pesquisa, e os conjuntos  $T_r$ ,  $A_r$  e  $U_r$  dos rótulos, atendentes, e usuários, respectivamente, que estão conectados a  $r$ , ou seja:  $T_r := \{t \mid (a, t, r, u) \in Y'\}$ ,  $A_r := \{a \mid (a, t, r, u) \in Y'\}$  e  $U_r := \{u \mid (a, t, r, u) \in Y'\}$ . Além disso, seja  $\text{tau}_r := \{(t, a, u) \mid (t, a, u) \in T \times A \times U \wedge (a, t, r, u) \in Y'\}$  a tripla (rótulo, atendente, usuário) ocorrendo com  $r$ . O coeficiente de agrupamento  $\gamma$  de  $r$  pode ser definido pela Equação (6.2):

$$\gamma(r) = \frac{|\text{tau}_r|}{|T_r| \cdot |A_r| \cdot |U_r|} \in [0, 1] \quad (6.2)$$

Sendo o numerador da Equação (6.2) o número de triplas (rótulos, atendentes, usuários) que ocorrem juntamente com o recurso  $r$ , e o denominador o número total de triplas (rótulos, atendentes, usuários), que poderiam ocorrer entre os rótulos, atendentes e usuários

que estão conectados a  $r$ . Se a vizinhança de  $r$  está maximamente agrupada, então todas as triplas de  $T_r \times A_r \times U_r$  vão ocorrer em  $tau_r$ , significando que o coeficiente de agrupamento é máximo (igual a 1 (um)).

A definição de coeficiente de agrupamento que foi mostrada acima, para um nó do tipo recurso da folksonomia obtida nesta pesquisa, vale simetricamente para atendentes (Equação (6.3)), rótulos (Equação (6.4)), e usuários (Equação (6.5)).

$$\gamma(a) = \frac{|tru_a|}{|T_a| \cdot |R_a| \cdot |U_a|} \in [0, 1] \quad (6.3)$$

$$\gamma(t) = \frac{|aru_t|}{|A_t| \cdot |R_t| \cdot |U_t|} \in [0, 1] \quad (6.4)$$

$$\gamma(u) = \frac{|atr_u|}{|A_u| \cdot |T_u| \cdot |R_u|} \in [0, 1] \quad (6.5)$$

O cálculo do coeficiente de agrupamento para todo o grafo  $G_F$  da folksonomia desta pesquisa, pode ser definido pela média das médias da Equação (6.6):

$$\gamma(G_F) = \text{média}\{\text{média}\{\gamma(a)\} + \text{média}\{\gamma(t)\} + \text{média}\{\gamma(r)\} + \text{média}\{\gamma(u)\}\}, \quad (6.6)$$

$$\forall a \in A, \forall t \in T, \forall r \in R, \text{ e } \forall u \in U$$

Após a obtenção do comprimento de caminho característico e do coeficiente de agrupamento da folksonomia aprendida pelo método, o próximo passo da avaliação de característica é comparar os valores dessas propriedades com as de um grafo gerado de modo aleatório. Para a geração do grafo aleatório é realizada uma adaptação da estratégia “*Binomial*” descrita em (CATTUTO et al., 2007), a qual originalmente foi desenvolvida para folksonomias do modelo tripartido. O intuito dessa adaptação é tornar a estratégia “*Binomial*” apta ao modelo quadripartido proposto nesta pesquisa. A geração de um grafo aleatório através da adaptação da estratégia “*Binomial*”, para o âmbito desta pesquisa, é realizada do seguinte modo:

- Os nós do grafo aleatório são obtidos a partir da extração de todos os nós dos conjuntos  $A$ ,  $T$ ,  $R$ ,  $U$ , de Atendentes, Rótulos, Recursos, e de Usuários, respectivamente, da folksonomia que foi gerada pelo método proposto. Na sequência

cria-se para o grafo aleatório a mesma quantidade de arestas que o grafo original gerado pelo método proposto possui. Para cada uma dessas arestas criadas, os nós em suas extremidades são selecionados através de distribuições uniformes, a partir dos nós dos conjuntos  $A$ ,  $T$ ,  $R$ ,  $U$  (extraídos para o grafo aleatório). O resultado disso é uma distribuição binomial sobre graus dos nós do grafo aleatório (CATTUTO et al., 2007).

Caso a folksonomia obtida pelo método desta pesquisa apresente um comprimento de caminho característico próximo e um coeficiente agrupamento grande, em relação a essas mesmas propriedades do grafo aleatório, pode-se concluir que a folksonomia desta pesquisa possui a característica de mundo-pequeno, conforme descrito no Capítulo 2.

Uma questão que deve ser destacada é que o cálculo do comprimento de caminho característico, para os grafos da folksonomia (e conseqüentemente aleatório) é um problema NP-completo, devido ao tamanho dos grafos. Assim, seguindo a abordagem de (CATTUTO et al., 2007), a propriedade de comprimento de caminho característico é calculada em relação a uma amostra aleatória de duzentos nós para cada um dos grafos gerados (folksonomia e aleatório). Além disso, vale ressaltar que a folksonomia desta pesquisa é representada por um modelo quadripartido  $\mathbb{F} := (A, T, R, U, Y')$ , deste modo o comprimento de caminho característico é obtido calculando-se as distâncias entre todos os tipos de nós do grafo (atendentes, rótulos, recursos e usuários). Com isso, além de se manter a estrutura original do modelo proposto e da folksonomia em si, também é mantida a integridade do relacionamento  $Y'$ .

### 6.3.2. Avaliação de Teste de Domínio

Para medir a acurácia da folksonomia aprendida pelo método, no âmbito da interpretação de enunciados de diálogos (verificando se eles pertencem ou não ao domínio representado pela folksonomia) é utilizado o método *Holdout* (TAN et al., 2005). O método *Holdout* foi aplicado dividindo-se o conjunto de dados em 2/3 para treinamento e 1/3 para teste do modelo. O conjunto de dados é representado por todos os diálogos do corpus de diálogos, sendo que 2/3 desses diálogos são utilizados para o aprendizado da folksonomia (são o corpus de entrada do método proposto) e os 1/3 restantes servem para medir a acurácia dela no processo de interpretação de enunciados.

O processo de interpretação de enunciados de diálogos pela folksonomia aprendida nesta pesquisa tem como intuito verificar se eles pertencem ou não ao domínio representado pela ela. A interpretação de enunciados ocorre de modo similar ao qual é realizada a Etapa de Obter Recursos da folksonomia, descrita na atividade de Aprender do método (Capítulo 5). Neste caso, também a partir de um dado enunciado de diálogo extraem-se todos os seus substantivos, se ele possuir mais do que 1 (um) conforme justificado na etapa de Obter Recursos da Folksonomia (Capítulo 5). Em seguida são verificados quantos desses substantivos são rótulos da folksonomia proposta. A obtenção do número de substantivos de um dado enunciado que são rótulos da folksonomia pode ser feita pela Taxa de Inclusão, calculada pela Equação (5.2) (Capítulo 5). Porém, no contexto desta avaliação de teste de domínio, os enunciados que tiverem suas Taxas de Inclusão  $p_{enu}$ , maior ou igual a uma Taxa de Inclusão  $p_1$ , definida pela Equação (5.3) do Capítulo 5, são considerados como pertencentes ao domínio representado pela folksonomia obtida com o método FolksDialogue. Conforme já descrito ao longo deste documento, nesta pesquisa existe a suposição de que os interlocutores do tipo atendente possuem conhecimento pleno do domínio. Assim, seus enunciados são considerados naturalmente como pertencentes ao domínio. Por outro lado, os interlocutores do tipo usuário buscam auxílio para a solução de uma tarefa, podendo requisitar qualquer coisa, mesmo que ela esteja fora do domínio compreendido pelos atendentes. Deste modo, nos 1/3 dos diálogos utilizados para o processo de interpretação, apenas os enunciados dos interlocutores do tipo usuário são adotados para a interpretação. O Quadro 14 exemplifica o processo de interpretação de enunciados de diálogos pela folksonomia aprendida pelo método proposto.

**Quadro 14** – Exemplo do processo de interpretação de enunciados de diálogos pela folksonomia aprendida.

Enunciados de Entrada	Substantivos Extraídos	Conjunto T de rótulos da folksonomia:	Cálculo da Taxa de Inclusão ( $p_{enu}$ )	Teste no domínio da folksonomia	Resultado
Olá, como vai?	{}, conjunto vazio.	$T := \{\text{férias, motivo, doença, aposentadoria, direito, função, remuneração, cargo, disposição, tempo}\}$	(Número de substantivos é zero)	<b>para um <math>p_1</math> igual a 75,00</b>	Fora do domínio (não possui nenhum substantivo).
O livro “Uma Breve História do Tempo” é interessante.	{livro, história, tempo}		$(1/3) = 33,33$		Fora do domínio (a Taxa de inclusão é menor que $p_1$ ).
Um servidor tem direito a alguma remuneração especial quando entra em férias?	{servidor, direito, remuneração, férias}		$(3/4) = 75,00$		Pertencente ao domínio (a Taxa de inclusão é maior ou igual a $p_1$ ).

Depois de calcular a Taxa de Inclusão para cada enunciado, dizendo se eles pertencem ou não ao domínio, cada um desses resultados é comparado com um dado rótulo atribuído por um especialista de domínio a cada um dos enunciados. A comparação é realizada através do cálculo da acurácia da Equação (6.7). Sendo nessa equação “*nCorreto*” o número de enunciados que tiveram o resultado de suas Taxas de Inclusões igual aos rótulos atribuídos pelo especialista, e “*n*” o número total de enunciados que estão sendo testados.

$$Acurácia = \left( \frac{nCorreto}{n} \right) \times 100 \quad (6.7)$$

Após a apresentação do processo de avaliação do método proposto, na próxima seção deste trabalho é descrita uma abordagem proposta para a detecção de tendências em folksonomias geradas pelo método FolksDialogue.

#### 6.4. Abordagem de Detecção de Tendências

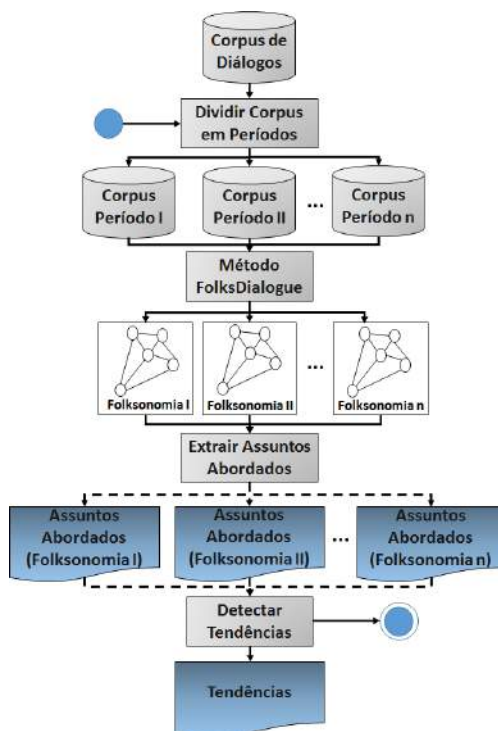
A detecção de tendências tem como um de seus intuitos verificar **quais são os assuntos abordados em diferentes intervalos de tempo**, no domínio dos diálogos do corpus de entrada do método FolksDialogue. Os assuntos abordados a serem detectados num dado intervalo de tempo são extraídos a partir de uma folksonomia aprendida com diálogos do corpus de entrada, pertencentes a apenas esse período de tempo. Assim, percebe-se que é construída uma folksonomia distinta para cada intervalo de tempo selecionado dentro do corpus de diálogos. Outro intuito da detecção de tendências desta pesquisa é confrontar cada um dos assuntos abordados, de uma dada folksonomia aprendida, com os assuntos das demais, verificando em quais folksonomias esses assuntos aparecem. O objetivo disso é devido a adoção da ideia de que se um dado assunto abordado aparece mais de uma vez ao longo do tempo (em folksonomias distintas) ele pode ser considerado como uma tendência. Ou seja, isso significa que os interlocutores do tipo usuário dos diálogos estão abordando tal assunto em diversos períodos de tempo. Além disso, outro ponto importante a se destacar é que cada um dos assuntos abordados extraídos de uma dada folksonomia aprendida está ranqueado (decrementemente), de acordo com o número de diálogos em que ele apareceu naquele intervalo de tempo. Com isso, torna-se possível verificar quais são os assuntos mais abordados num dado período, sendo o primeiro assunto do ranque o mais abordado, o segundo sendo considerado o segundo mais abordado, e assim sucessivamente. Por fim, na

detecção de tendências desta pesquisa ainda é possível verificar se um dado assunto abordado teve perda ou ganho de popularidade ao longo de intervalos de tempo distintos. Isso pode ser realizado verificando se um determinado assunto abordado apareceu em diferentes folksonomias, e se ele mudou de posição nos ranques dessas folksonomias.

A Figura 18 mostra a abordagem proposta. O primeiro passo para a detecção dos assuntos e das possíveis tendências, que eles possam ser ao longo de um intervalo de tempo, é pegar o corpus de diálogos de um dado domínio e particioná-lo em “períodos” ou “intervalos de tempo”. Para o particionamento os diálogos devem possuir uma informação que identifique o período em que eles foram produzidos, ou apenas estar dispostos em ordem cronológica no corpus. O número de partições a serem feitas depende do período em que se deseja extrair os assuntos/tendências.

Após a divisão do corpus de diálogos em intervalos de tempo, cada partição resultante é utilizada como entrada do método *FolksDialogue* visando a criação de uma folksonomia distinta. Em seguida, a partir de cada uma das folksonomias aprendidas realiza-se a extração de seus assuntos abordados. O processo de extração dos assuntos abordados em uma dada folksonomia é realizado com base em seus conjuntos  $T$ ,  $U$  e  $R$  de rótulos, de usuários, e de recursos, respectivamente. Vale destacar que o conjunto  $U$  é resultado do uso de uma característica dos diálogos orientados a tarefas, que é o tipo de interlocutor usuário, o qual busca auxílio na solução de uma dada tarefa (descrito na seção 3.1). A seguir são apresentadas as definições necessárias para a compreensão do modo como são obtidos os assuntos abordados.





**Figura 18** – Diagrama simbolizando o fluxo da abordagem proposta para detecção de tendências.

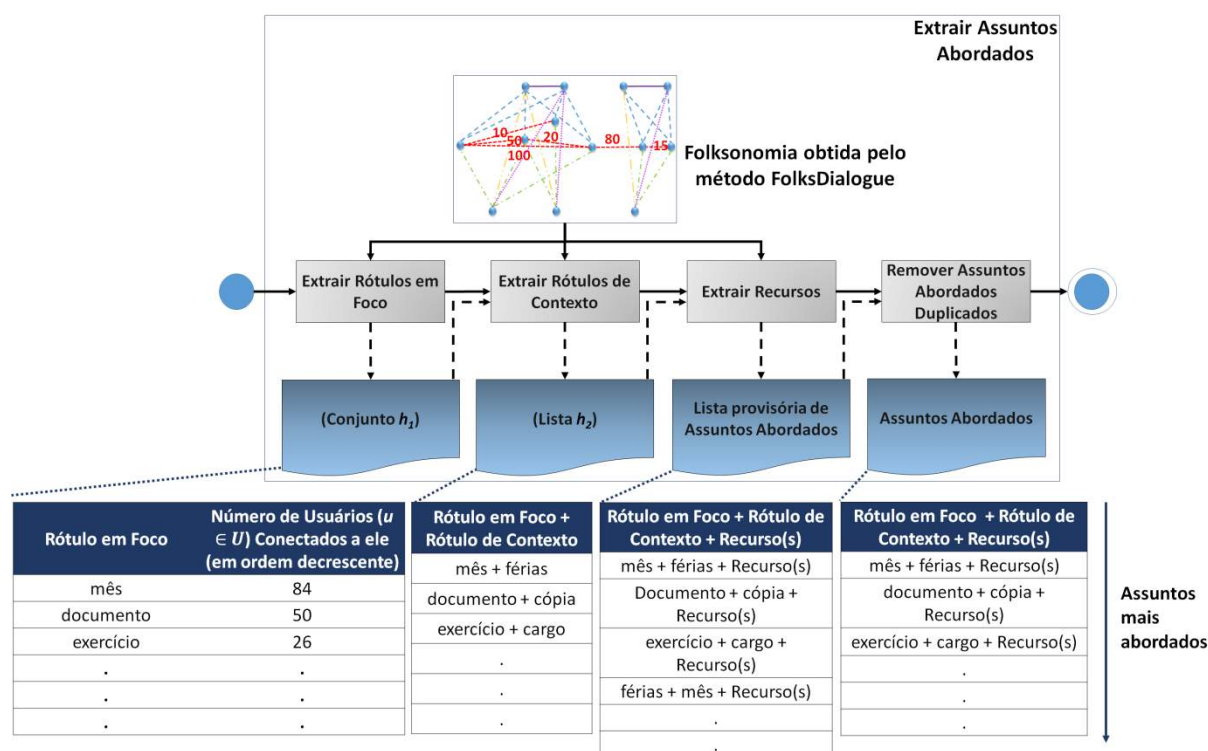
**Definição 6.1.** Um **Rótulo em Foco** é um rótulo  $t$  pertencente ao conjunto  $T$  da folksonomia  $\mathbb{F}$  descrita neste trabalho, que possui um número de usuários  $(u \in U)$  conectados a ele.

**Definição 6.2.** Um **Rótulo de Contexto** é um rótulo  $t$  pertencente ao conjunto  $T$  da folksonomia  $\mathbb{F}$  desta pesquisa, o qual está conectado a um dado Rótulo em Foco, através de um relacionamento  $b$  do conjunto  $B$  de relacionamentos dos rótulos de  $\mathbb{F}$ .

**Definição 6.3.** Um **Assunto Abordado** é composto por um Rótulo em Foco, por um Rótulo de Contexto, e por recursos  $(r \in R)$ , possuindo a seguinte nomenclatura: **Rótulo em Foco + Rótulo de Contexto + Recurso(s)**. Esses recursos que estão associados aos Rótulos em Foco e de Contexto são os recursos que ambos rotulam juntos. O objetivo do Rótulo de Contexto e dos recursos é contextualizar o Rótulo em Foco, formando assim um Assunto Abordado.

A Figura 19 ilustra em detalhes o processo realizado pelo bloco “Extrair Assuntos Abordados” da Figura 18. A obtenção dos Assuntos Abordados é realizada extraindo-se da folksonomia aprendida, um conjunto  $h_l$  com todos os rótulos  $(t \in T)$  que ela possui. Os rótulos de  $h_l$  estão ordenados decrescentemente com base no número de usuários  $(u \in U)$  que

estão conectados a eles. Isso significa que os rótulos que possuírem mais usuários conectados a eles estarão no topo ou no início do conjunto extraído. Nesta pesquisa é adotado que os interlocutores do tipo usuário, dos diálogos orientados a tarefas, são únicos, ou seja, cada diálogo é composto por um usuário distinto. Deste modo, pode-se inferir que o número de usuários ligados a um rótulo, representa o número de diálogos distintos em que tal rótulo apareceu, sendo os rótulos do início de  $h_1$  os mais utilizados em diálogos diferentes. O objetivo disso é preparar os assuntos abordados que estão sendo formados a estarem ranqueados, de modo que no topo do ranque estarão os assuntos que foram mais abordados.



**Figura 19** – Detalhes do processo que realiza a extração dos Assuntos Abordados.

Os rótulos de  $h_1$  são denominados de Rótulos em Foco (Definição 6.1). Após a aquisição dos Rótulos em Foco, o próximo passo para formar os Assuntos Abordados de uma dada folksonomia é obter o Rótulo de Contexto de cada um dos Rótulos em Foco de  $h_1$ . O motivo disso é que caso os Assuntos Abordados fossem formados apenas pelos Rótulos em Foco, eles poderiam não descrever com precisão os assuntos extraídos. Exemplificando essa questão, no âmbito de diálogos referentes ao domínio de Recursos Humanos, se um Assunto Abordado fosse formado apenas pelo Rótulo em Foco “mês”, esse assunto poderia ser referente a diversos temas, como mês de férias, mês de licença, mês de aposentadoria, dentre

outros. No entanto, não seria possível saber a qual desses temas o assunto estaria se referindo.

Assim, com o objetivo de melhorar a precisão dos assuntos que estão sendo detectados é realizada a obtenção dos Rótulos de Contexto. Dado um Rótulo em Foco de  $h_1$ , verifica-se na folksonomia aprendida quais são os rótulos ( $t \in T$ ) que possuem um relacionamento ( $b \in B$ ) com este Rótulo em Foco. O rótulo  $t$  da folksonomia que possuir o relacionamento de maior peso com esse Rótulo em Foco será o Rótulo de Contexto dele. A medida que os Rótulos em Foco tiverem seus respectivos Rótulos de Contexto encontrados, eles são armazenados numa lista  $h_2$ . Caso dois ou mais rótulos de  $T$  possuam o mesmo peso de ligação com um dado Rótulo em Foco, todos eles serão Rótulos de Contexto e cada um estará formando um Assunto Abordado distinto com esse Rótulo em Foco. Os Rótulos em Foco que não possuírem nenhum relacionamento com outros rótulos ( $t \in T$ ) da folksonomia aprendida são descartados de  $h_1$ , e conseqüentemente não formarão Assuntos Abordados. A razão para isso é a suposição de que um Rótulo em Foco, que apareceu numa larga gama de diálogos distintos tem uma grande probabilidade de possuir relacionamentos com outros rótulos da folksonomia aprendida. Por outro lado, acredita-se que Rótulos em Foco que não possuem relacionamentos com outros rótulos de  $R$ , tendem a aparecer numa quantidade menor de diálogos, não estando deste modo entre os assuntos mais abordados pelos interlocutores.

Com a lista  $h_2$  gerada, o próximo passo é pegar cada elemento dela, ou seja, um dado “par” de Rótulo em Foco e de Rótulo de Contexto, e buscar os recursos ( $r \in R$ ) que ambos esses rótulos rotulam juntos. Conforme descrito na Definição 6.3, os recursos, assim como os Rótulos de Contexto, têm como intuito contextualizar um dado Rótulo em Foco. Após a obtenção dos recursos de cada “par” de Rótulo em Foco e Rótulo de Contexto, os Assuntos Abordados são considerados formados, porém neste passo atual gera-se uma lista denominada Lista Provisória de Assuntos Abordados. O passo final para a obtenção dos Assuntos Abordados é pegar a Lista Provisória de Assuntos Abordados e verificar se existem Assuntos Abordados que estão duplicados nela. Ou seja, dados dois determinados Assuntos Abordados  $A_1$  e  $A_2$ , diz-se que eles estão duplicados, se o Rótulo em Foco de  $A_1$  for igual ao Rótulo de Contexto de  $A_2$ , e se o Rótulo em Foco de  $A_2$  for igual ao Rótulo de Contexto de  $A_1$ , ou vice-versa, isto é,  $A_1 = \text{mês} + \text{férias} + \text{recurso(s)}$  e  $A_2 = \text{férias} + \text{mês} + \text{recurso(s)}$ . Caso sejam encontrados dois Assuntos Abordados duplicados, um deles é escolhido de forma arbitrária e removido de lista provisória. Após a remoção dos possíveis assuntos duplicados, o resultado é uma lista contendo os Assuntos Abordados obtidos pela abordagem de detecção de tendências

desta pesquisa. A lista de Assuntos Abordados está disposta em ordem decrescente, isto é, no topo dela estão os assuntos que foram mais abordados pelos interlocutores dos diálogos do corpus de entrada do método proposto.

Após a construção das folksonomias com diálogos pertencentes a diferentes intervalos de tempo, e da extração dos assuntos abordados a partir de cada uma delas, um especialista de domínio (do corpus de diálogos) avalia se os assuntos extraídos realmente representam as interlocuções do domínio. Para isso, o especialista confronta os rótulos (em Foco e de Contexto) com os recursos (enunciados) de cada um dos Assuntos Abordados, e verifica se o Assunto Abordado reflete o conteúdo desses enunciados.

Uma vez obtidos os Assuntos Abordados, a última parte da abordagem de detecção de tendências é verificar quais deles podem ser possíveis tendências ao longo de um intervalo de tempo. Para isso, cada um dos Assuntos Abordados de uma dada folksonomia aprendida é confrontado com os Assuntos Abordados das demais, verificando-se assim em quais folksonomias ele aparece. Se um dado Assunto Abordado aparecer mais de uma vez ao longo do tempo (em folksonomias distintas) ele pode ser considerado como uma tendência.

## **6.6. Conclusão**

Neste capítulo foi apresentado o processo para implementação das atividades e suas respectivas etapas que compõem o método proposto. O processo de implementação é realizado através de um protótipo computacional, composto pelas atividades e etapas do método, e tem como intuito validar o FolksDialogue. Além disso, foram apresentadas duas avaliações para o método proposto, uma denominada avaliação de característica e outra de avaliação de teste de domínio. Na sequência, foi proposto um modo de se detectar e extrair tendências das folksonomias geradas pelo método.

No próximo capítulo são mostrados os resultados obtidos com o método FolksDialogue.

# Capítulo 7

## Resultados Obtidos

Neste capítulo são descritos os resultados obtidos por esta pesquisa. Os resultados foram obtidos através da realização de três experimentos. O primeiro experimento visa implementar a avaliação de característica descrita no Capítulo 6: verificar que as estruturas geradas pelo método proposto possuem o aspecto de mundo-pequeno. Um segundo experimento foi efetuado para implementar a avaliação de teste de domínio do Capítulo 6. Essa avaliação visa mostrar que as folksonomias obtidas pelo FolksDialogue podem ser usadas para se interpretar enunciados de diálogos, verificando se eles pertencem ou não ao domínio que elas representam. Por fim, foi elaborado um terceiro experimento para aplicar a abordagem de detecção de tendências proposta por esta pesquisa. A apresentação dos resultados é realizada do seguinte modo: inicialmente é descrito o corpus de entrada utilizado nos experimentos (seção 7.1), e na sequência são mostrados os resultados produzidos pelos experimentos de avaliação de característica (seção 7.2), experimento de avaliação de teste de domínio (seção 7.3) do método proposto, e resultados da abordagem de detecção de tendências (seção 7.4).

### 7.1. Corpus de Entrada

Para os experimentos realizados foram selecionados os 901 diálogos do corpus da prefeitura descrito na seção 6.1. Vale ressaltar que os diálogos foram produzidos durante quatro anos (2006 a 2009), e estão dispostos em ordem cronológica (seção 6.1). Nestes diálogos são tratados assuntos como aposentadoria, direitos de ordem geral, estágio

probatório, férias, licenças e provimento. Os 901 diálogos do corpus são compostos de um total de 7.064 enunciados. Os diálogos do corpus foram produzidos por 5 (cinco) interlocutores do tipo atendente, e cada diálogo teve um interlocutor do tipo usuário distinto, totalizando 901 usuários. No Quadro 15 são apresentados exemplos destes diálogos.

**Quadro 15** – Exemplos de Diálogos Utilizados como Corpus de Entrada no Método Proposto.

ID do Diálogo	Tipo de Interlocutor	Enunciado
1	$u_1$	Bom dia.
1	$a_1$	Bom dia, em que posso ajudar?
1	$u_1$	Quais doenças são consideradas graves, contagiosas ou incuráveis?
1	$a_1$	hanseníase, cardiopatia grave, doença de Parkinson, paralisia irreversível e incapacitante, espondiloartrose anquilosante,
1	$u_1$	nefropatia grave, estados avançados do mal de Paget (osteíte deformante), AIDS, e outras que e lei indicar, com base na medicina especializada.
1	$a_1$	Ok.
1	$u_1$	Se tiver mais alguma dúvida, fique à vontade para perguntar.
2	$u_2$	Boa tarde, faço 70 anos o ano que vem, tenho que me aposentar?
2	$a_2$	sim , sua aposentadoria será compulsoria.
2	$u_2$	E eu preciso fazer algum pedido para isso?
2	$a_2$	O estatuto diz que a aposentadoria compulsória será automática, e declarada por ato, com vigência a partir do dia imediato àquele em que o servidor atingir a idade limite de permanência no serviço ativo.
2	$u_2$	Obrigada.
2	$a_2$	Estamos à disposição.
3	$u_3$	Boa tarde. Eu quero saber uma coisa.
3	$a_1$	Pode perguntar.
3	$u_3$	Todo servidor tem direito a licença p/ tratar de assuntos particulares?
3	$a_1$	De acordo com o Estatuto, não se concederá esta licença ao servidor que esteja respondendo à sindicância, processo administrativo ou, a qualquer título, esteja ainda obrigado à indenização ou à devolução aos cofres públicos.
3	$u_3$	Obrigada por responder
3	$a_1$	Por nada.
4	$u_4$	Oi!
4	$a_3$	Olá. Em que posso ajudar?
4	$u_4$	Como faço para receber os dez dias de férias?
4	$a_3$	Basta apresentar um requerimento ao recursos humanos, desde que este seja apresentado 30 dias antes do período de fruição do período aquisitivo correspondente aos dias vendidos.
4	$u_4$	Obrigada
4	$a_3$	Estamos à disposição.

## 7.2. Resultados da Avaliação de Característica

Para verificar se a folksonomia gerada pelo método FolksDialogue possui o aspecto de mundo-pequeno, foram calculados seu comprimento de caminho característico (Equação (6.1) deste documento) e seu coeficiente de agrupamento (Equação (6.6) deste documento). Os resultados destes cálculos foram comparados com os de um grafo gerado aleatoriamente a partir de uma distribuição binomial (descrito na seção (6.3)).

Primeiramente foi realizado o aprendizado de uma folksonomia através método FolksDialogue. A partir da aplicação do corpus de diálogos descrito na seção 7.1, foi gerada uma folksonomia resultante pelos parâmetros apresentados na Tabela 1 com 1.658 nós. Sendo que destes, 5 eram atendentes, 289 eram rótulos, 463 eram recursos, e 901 eram usuários. Na sequência foram extraídos os nós da folksonomia, e partir deles construiu-se um grafo aleatório baseado em distribuição binomial (descrito na seção 6.3).

**Tabela 1** - Parâmetros da Folksonomia Aprendida para a Avaliação de Característica.

Parâmetro	Valor
Nº de Substantivos extraídos dos Enunciados dos Atendentes	12.408
Conjunto de Substantivos ( $L_s$ )	560
$fc_1$	4,94
$fc_2$	2,32
$p_1$	3,30
Nº de Relacionamentos entre os Rótulos	140

A Tabela 2 mostra os resultados dos comprimentos de caminho característico e dos coeficientes de agrupamento para a folksonomia e para o grafo aleatório. A partir desses resultados é possível concluir que a folksonomia possui o aspecto de mundo-pequeno, isto é, seu coeficiente de agrupamento é maior, e seu comprimento de caminho característico é comparável aos de um grafo aleatório. O motivo para o coeficiente de agrupamento do grafo aleatório ter sido 0,0 pode ser atribuído ao fato da folksonomia ter apenas 1.658 nós, e assim torna-se improvável que ele forme exatamente as triplas dos numeradores das Equações (6.3), (6.4), (6.5), e (6.6). Como exemplo de comparação, em (CATTUTO et al., 2007) o coeficiente de agrupamento do grafo aleatório foi de aproximadamente 0,2 para um grafo com  $2^{17}$  atribuições realizadas pelos usuários.

**Tabela 2** – Resultados da avaliação de característica.

	<b>Folksonomia</b>	<b>Grafo Aleatório</b>
Comprimento de Caminho Característico	2,9	2,7
Coefficiente de Agrupamento	0,5	0,0

Uma das conclusões que pode ser extraída em relação ao comprimento de caminho característico pequeno, numa folksonomia obtida a partir de diálogos, é que o conteúdo dos diálogos (utilizados para o aprendizado) está conectado entre si. Isso significa que os conteúdos tem correlação entre si, e modelam um domínio de forma convergente e consistente. No caso do coeficiente de agrupamento ser grande numa folksonomia aprendida com diálogos, a ideia pode ser análoga a descrita na seção 2.2, porém neste caso conclui-se que os interlocutores dos diálogos utilizaram o mesmo vocabulário em seus enunciados, fazendo-o de forma coerente e consistente.

### 7.3. Resultados da Avaliação de Teste de Domínio

O primeiro passo para medir a acurácia da folksonomia, na interpretação de enunciados de diálogos (verificando se eles pertencem ou não ao domínio que ela representa), foi dividir o corpus de entrada através da abordagem *Holdout*. Dos 901 diálogos do corpus, 598 (2/3 do total, extraídos em ordem cronológica) foram destinados para o aprendizado da folksonomia, e os 303 restantes (1/3 do corpus) foram usados para medir a acurácia dela no processo de interpretação de enunciados.

O uso dos 598 diálogos como entrada do método *FolksDialogue* resultou numa folksonomia composta por cinco personomias. Durante a execução do método foram obtidos 2.401 enunciados dos atendentes (50,97% de todos os enunciados). Destes enunciados foram extraídos 542 substantivos únicos, os quais formaram o conjunto *Ls*. Usando uma frequência de corte  $fc_1 = 4,83$ , de todos os substantivos de *Ls*, 278 (51,29%) viraram rótulos da folksonomia. Usando uma Taxa de Inclusão  $p_1 = 3,34$ , de todos os enunciados produzidos pelos atendentes, 155 (6,46%) se tornaram recursos da folksonomia. A respeito do relacionamento entre rótulos, a partir de uma frequência de corte  $fc_2 = 1,21$ , o método criou 41 relacionamentos entre rótulos.

Após realizar o aprendizado da folksonomia, a acurácia dela no processo de interpretação de enunciados de diálogos foi medida. Dos 2.354 enunciados que eram parte dos



303 diálogos usados nesse processo de interpretação, extraíram-se todos os 1156 (49,11%) enunciados dos interlocutores do tipo usuário. Para cada um desses 1156 enunciados foram calculadas suas Taxas de Inclusão Equação (5.2) (Capítulo 5), e em seguida aplicada a Taxa de Inclusão  $p_1$  (a mesma gerada pela aprendizagem), determinando quais enunciados faziam parte do domínio em questão. Na sequência, esses resultados foram comparados com os rótulos atribuídos por um especialista, e então aplicados a acurácia da Equação (6.7), conforme descrito na seção 6.3. A aplicação dos 1156 enunciados no processo de interpretação de enunciados produziu um acurácia de 69,20%. Os 30,88% que representam o erro de interpretação são compostos por 29,23% de falsos negativos e 1,65% de falsos positivos. A explicação para que a grande maioria do erro tenha sido de falsos negativos pode ser devido à divisão do corpus em treinamento e teste. Dado a esse particionamento, parte do domínio dos enunciados de teste pode não estar na partição de treinamento, e consequentemente eles são interpretados como fora do domínio. No entanto, vale ressaltar que o modelo final da folksonomia construída (não na ótica da avaliação) adota 100% do corpus de diálogos.

Conforme mencionado no Capítulo 1, em comparação com outras estruturas (como ontologias), folksonomias são formas de representação do conhecimento mais simples de se implementar e utilizar. Assim o resultado da acurácia pode ser afetado pelo fato de que folksonomias possuem somente substantivos (seus rótulos), os quais são usados no processo de interpretação de enunciados de diálogos. Isto é, os substantivos podem não conter toda a semântica de um enunciado. Outro fator que pode impactar na acurácia é a maneira como os interlocutores do tipo usuário escrevem seus enunciados. Em muitos casos eles usam abreviações (“p”, ao invés de “para”), erros de ortografia (“qauntos”, ao invés de “quantos”), siglas (“DGP”), e gírias (“Vlw”, que significa “Valeu”), ocasionando em problemas para o parser. Além disso, as partições do corpus de diálogos usadas para o aprendizado e para interpretar enunciados, podem não garantir que todo o domínio está presente em ambas delas. Isto é, os diálogos usados para o aprendizado da folksonomia podem ser relacionados a determinados assuntos e conter termos característicos. Por outro lado, parte da partição usada para a interpretação de enunciados pode ser peculiar a outros assuntos (que não fazem parte dos diálogos da partição de aprendizagem) com seus termos específicos. Consequentemente, durante o processo de interpretação de enunciados, os termos dos enunciados que estão sendo interpretados podem não ser encontrados na folksonomia aprendida, e a acurácia pode

diminuir. Uma solução para este viés poderia ser o uso da validação cruzada (TAN et al., 2005). No entanto, devido a sua característica de divisão do conjunto de dados em diversos subconjuntos, e normalmente se utilizando extratificação, para garantir que cada classe dos exemplos seja representada em proporções iguais em todos os subconjuntos, isto elevaria a um alto custo computacional. Além disso, o processo de extratificação poderia quebrar a estrutura dos diálogos (atendentes, usuários, enunciados, etc.), o que afetaria a aprendizagem da folksonomia pelo método FolksDialogue. Um exemplo disso, é que as entidades da folksonomia devem estar devidamente ligadas entre si. Isto é, no caso dos atendentes, partir dos diálogos do corpus, eles devem estar conectados aos seus enunciados e aos usuários com quem dialogaram. Deste modo, optou-se pela não utilização da validação cruzada.

Para confirmar a hipótese descrita acima de que o domínio em questão pode não estar presente simultaneamente em ambas as partições (aprendizado e teste), e conseqüentemente afetar a acurácia, foi realizado um experimento que alterou a proporção da divisão de 2/3 (aprendizagem) e 1/3 (teste) da abordagem *holdout*. Neste experimento foi arbitrado que a divisão do corpus seria de 90% de seu total (extraídos cronologicamente) para o aprendizado da folksonomia e os 10% restantes para o teste de domínio dos enunciados. Assim, o modelo aprendido, ou seja, a folksonomia do método FolksDialogue ficaria mais robusto e consistente. Dos 901 diálogos do corpus, 811 (90%) foram utilizados para o aprendizado da folksonomia, e os 90 restantes (10%) foram adotados para medir a acurácia dela no processo de interpretação de enunciados. O uso dos 811 diálogos como entrada do método FolksDialogue resultou novamente numa folksonomia composta por cinco personomias. Durante a execução do método foram obtidos 3.129 enunciados dos atendentes (49,04% de todos os enunciados). Destes enunciados foram extraídos 551 substantivos únicos, os quais formaram *Ls*. Usando uma frequência de corte  $fc_1 = 4,90$ , de todos os substantivos de *Ls*, 294 (53,36%) viraram rótulos da folksonomia. Usando uma Taxa de Inclusão  $p_1 = 3,62$ , de todos os enunciados produzidos pelos atendentes, 449 (12,35%) se tornaram recursos da folksonomia. A respeito do relacionamento entre rótulos, a partir de uma frequência de corte  $fc_2 = 2,31$ , o método criou 138 relacionamentos entre rótulos.

Em relação ao processo de interpretação de enunciados, dos 684 enunciados que eram parte dos 90 diálogos usados nesse processo de interpretação, extraíram-se todos os 336 (49,12%) enunciados dos interlocutores do tipo usuário. Para cada um desses 336 enunciados foram calculadas suas Taxas de Inclusão Equação (5.2) (Capítulo 5), e em seguida aplicada a

Taxa de Inclusão  $p_1$  (a mesma gerada pela aprendizagem), determinando quais enunciados faziam parte do domínio em questão. Na sequência, esses resultados foram comparados com os rótulos atribuídos por um especialista, e então aplicados a acurácia da Equação (6.7), conforme descrito na seção 6.3 deste documento. A aplicação dos 336 enunciados no processo de interpretação de enunciados produziu um acurácia de 72,32%. Os 27,68% que representam o erro de interpretação são compostos por 26,19% de falsos negativos e 1,49% de falsos positivos. Com isso foi possível comprovar que se ambas as partições extraídas do corpus (aprendizado e teste) não possuem o mesmo domínio, a acurácia da interpretação de enunciados de diálogos pode ser afetada.

#### 7.4. Resultados da Abordagem de Detecção de Tendências

Nesta seção são apresentados os resultados obtidos com a implementação da abordagem de detecção de tendências proposta no Capítulo 6. Inicialmente o corpus de diálogos descrito na seção 7.1 foi dividido em “intervalos de tempo”. O corpus foi dividido em oito partes ou oito “intervalos de tempo”, cada um representando seis meses do corpus (visto que o período completo é quatro anos). A escolha por esse número de partições foi para que cada intervalo de tempo representasse seis meses. Cada uma das oito partes foi usada como entrada do método FolksDialogue, resultando em oito folksonomias distintas (I, II, III, IV, V, VI, VII, VIII). Na sequência, as oito folksonomias foram aplicadas na abordagem de detecção de tendências proposta por esta pesquisa (seção 6.4), e seus Assuntos Abordados foram extraídos (Tabela 3). Cada um dos assuntos abordados foi validado por um especialista de domínio, que confrontou os rótulos (em Foco e de Contexto) com os recursos (enunciados) de cada um dos Assuntos Abordados, e verificou que eles refletiam o conteúdo desses enunciados.

**Tabela 3** – Assuntos Abordados extraídos por cada folksonomia.

Folksonomia	Nº de Assuntos Abordados
I	11
II	1
III	14
IV	34
V	39
VI	67
VII	49
VIII	63

O número de Assuntos Abordados de cada folksonomia distingue devido ao fato de que elas foram aprendidas com diálogos pertencentes a diferentes intervalos de tempo. Em cada período os interlocutores do tipo atendente, que produziram os enunciados, podem ter sido diferentes, e conseqüentemente, a maneira como eles geram seus enunciados também distingue.

Um dos motivos da “Folksonomia II” ter apenas um Assunto Abordado pode ser creditado em particular ao modo como os interlocutores do tipo atendente escreveram seus enunciados naquele intervalo de tempo. Por exemplo, comparando a “Folksonomia I” com a “Folksonomia II”, a primeira é composta por 192 rótulos e 106 recursos, enquanto que a segunda tem 168 rótulos e apenas 64 recursos. A razão para isso é que nos enunciados utilizados para o aprendizado da “Folksonomia II”, os termos usados pelos interlocutores do tipo atendente não foram considerados importantes pelo IDF (etapa de “Obter Rótulos da Folksonomia” do método, Capítulo 5). Quando uma folksonomia tem menos rótulos, a probabilidade de um enunciado se tornar um recurso é menor, e conseqüentemente a probabilidade de existir um relacionamento  $b \in B$  entre dois dados rótulos também é menor. Assim, formam-se menos Assuntos Abordados devido a não existência de ligação entre Rótulos em Foco e Rótulos de Contexto.

Após extrair os Assuntos Abordados de todos os intervalos de tempo, a abordagem da detecção de tendências procurou por possíveis tendências nesses intervalos, isto é, se um dado Assunto Abordado apareceu em diferentes intervalos de tempo. Foram encontrados 39 Assuntos Abordados que se repetiram ao longo do tempo. O Quadro 16 mostra alguns dos Assuntos Abordados que se tornaram tendências. Exemplificando, o Assunto Abordado “Registro + Problema + recurso(s)” apareceu em três intervalos de tempo (representados pelas folksonomias I, III, VII). Isto significa que nos diálogos do corpus entrada, os interlocutores do tipo usuário reportaram problemas de registro nos primeiros seis meses de 2006 (folksonomia I), 2007 (folksonomia III), e 2009 (folksonomia VII). A existência dessa tendência poderia ser útil para alertar alguém do domínio sobre tal recorrente problema.

**Quadro 16** – Exemplo de Assuntos Abordados que se tornaram tendências.

<b>Tendência (Assunto Abordado)</b>	<b>Folksonomias Contendo a Tendência</b>
Registro + Problema + recurso(s)	I, III, VII
Filho + Nascimento+ recurso(s)	IV, VIII
Classificação + Carreira + recurso(s)	III, VII
Teste + Aplicação + recurso(s)	VI, VII

## **7.5. Conclusão**

Este capítulo apresentou os resultados obtidos por esta pesquisa a partir da implementação do método FolksDialogue. Foram apresentados os resultados gerados por três experimentos distintos: o primeiro implementou a avaliação de característica, o segundo a avaliação de teste de domínio, e o último a abordagem de detecção de tendências desta pesquisa.

# Capítulo 8

## Conclusões e Trabalhos Futuros

Sistemas de diálogo visam tornar mais natural a interação dos seres humanos com os computadores. Em sistemas de diálogos um dos componentes fundamentais do processo de interpretação dos enunciados, apresentados pelos usuários, é o modelo conceitual. O modelo conceitual descreve um dado domínio e pode ser especificado através de diversas formas de representação do conhecimento. Nesta pesquisa foi proposto representar o modelo conceitual através de folksonomias.

A obtenção das folksonomias a serem utilizadas para especificar o modelo conceitual é realizada pelo método descrito nesta pesquisa, o qual visa efetuar o aprendizado automático de folksonomias a partir de diálogos orientados à tarefa em português do Brasil. As folksonomias construídas através de diálogos são representadas por um modelo quadripartido proposto neste trabalho, constituído de Atendentes, Rótulos, Recursos, e Subconjuntos de Usuários. O método que realiza o aprendizado das folksonomias é composto pelas atividades de Pré-processar e de Aprender, as quais são realizadas de forma automática. Computacionalmente, as folksonomias geradas pelo método proposto são representadas através de grafos.

Com o intuito de validar e avaliar o método proposto por esta pesquisa foram definidas duas formas de avaliações, uma denominada de avaliação de característica, e a outra de avaliação de teste de domínio. A avaliação de característica teve por objetivo verificar se as folksonomias geradas pelo método FolksDialogue possuíam o aspecto de mundo-pequeno, o qual é uma característica das folksonomias (para alguns autores da literatura, descrito no Capítulo 6). Por outro lado, a avaliação de teste de domínio visou medir a acurácia da folksonomia aprendida pelo método, no âmbito da interpretação de enunciados de diálogos,

verificando se eles pertencem ou não ao domínio representado pela folksonomia.

Além disso, nesta pesquisa foi apresentada uma abordagem para realizar a detecção de tendências nas folksonomias geradas pelo método proposto. A detecção de tendências objetivou verificar quais eram os assuntos abordados em diferentes intervalos de tempo, no domínio dos diálogos do corpus de entrada do método FolksDialogue.

Foram realizados três experimentos os quais produziram os resultados obtidos por esta pesquisa. O primeiro experimento teve como intuito implementar a avaliação de característica descrita nesta pesquisa. O objetivo do segundo experimento foi verificar qual era a acurácia da folksonomia no processo de interpretação de enunciados de diálogos, através da implementação da avaliação de teste de domínio. E por fim, o terceiro experimento visou implementar a abordagem de detecção de tendências proposta por esta pesquisa.

Através do primeiro experimento, que implementou a avaliação de característica, foi possível comprovar que as estruturas geradas pelo método proposto possuem o aspecto de mundo-pequeno. A comprovação foi constatada pelo fato das folksonomias apresentarem um comprimento de caminho característico comparável, e possuírem um coeficiente de agrupamento considerado grande, ambos em relação aos de um grafo aleatório.

O segundo experimento, o qual implementou a avaliação de teste de domínio, mostrou que a folksonomia do método proposto possui uma acurácia de 72,32% no processo de interpretação de enunciados de diálogos. Um dos fatores que pode ter influenciado e afetado o resultado da acurácia pode ser devido as folksonomias serem formas de representação do conhecimento simples (se comparadas com outras estruturas, como ontologias), possuindo apenas termos ou palavras-chave, ou seja, seus rótulos, os quais são usados no processo de interpretação de enunciados de diálogos. Isto é, os substantivos podem não conter toda a semântica de um enunciado. Também, o modo como os interlocutores dos diálogos escrevem seus enunciados, produzindo-os com erros ortográficos, e fazendo uso de abreviações, pode ter acarretado em problemas para o parser, influenciado no resultado da acurácia no processo de interpretação. Além disso, as partições do corpus de diálogos usadas para o aprendizado da folksonomia, e para o processo de interpretação de enunciados, podem não ter garantido que todo o domínio estivesse presente em ambas delas. Cada partição pode ter em sua composição determinados termos que não estão presentes na outra, afetando desse modo a acurácia.

O terceiro experimento realizado implementou a abordagem de detecção de tendências proposta por esta pesquisa. A partir de folksonomias construídas com diálogos pertencentes a

diferentes intervalos de tempo, foi possível extrair os assuntos abordados pelos interlocutores nos diálogos utilizados para o aprendizado dessas estruturas. Os assuntos abordados que tiveram aparição ao longo de diferentes intervalos de tempo foram caracterizados como tendências.

Deste modo, com base no método FolksDialogue, nos experimentos e avaliações realizadas, foi possível se atingir os objetivos principal e específicos propostos por esta pesquisa. Além disso, as duas hipóteses apresentadas por este trabalho: i) “*É possível construir uma folksonomia a partir de diálogos.*” e ii) “*Com a folksonomia construída, é possível indicar se enunciados de diálogos pertencem ou não ao domínio representado por ela.*” puderam ser comprovadas. A validação da primeira ocorreu pela projeção e implementação do método FolksDialogue, que realiza a construção de folksonomias a partir de diálogos, já a segunda hipótese foi comprovada pelos resultados da avaliação de teste de domínio (Capítulo 7) que indica se enunciados pertencem ou não ao domínio da folksonomia.

Um dos trabalhos futuros desta pesquisa é o aprimoramento da forma como é realizado o processamento de linguagem natural do método proposto. O objetivo é corrigir questões relacionadas a erros ortográficos e abreviações nos enunciados dos diálogos. Este aprimoramento poderia servir de auxílio no processo de interpretação de enunciados de diálogos, pois os termos testados estariam padronizados de acordo com a gramática da língua portuguesa, consequentemente a acurácia da interpretação poderia ser aumentada. Também, espera-se testar o método proposto com outros corpora de diálogos em português do Brasil. O intuito é poder analisar como o método proposto se comportará e quais dificuldades poderão surgir. Além disso, existe o intuito de se verificar a possibilidade da realização do aprendizado de ontologias a partir de diálogos. Com isso, poderá ser possível estudar e comparar as vantagens e desvantagens das folksonomias e das ontologias no aprendizado de diálogos. No entanto, o aprendizado de ontologias a partir de diálogos pode ser um desafio, dado que ontologias possuem um elevado formalismo e são estruturas mais complexas que folksonomias. Por fim, no âmbito da abordagem de detecção de tendências um dos trabalhos futuros que pode ser realizado é o estudo de *concept drift* (TSYMBAL, 2004). Dado ao fato de que não existe uma garantia do comportamento dos usuários nos diálogos e consequentemente uma estabilidade nas tendências extraídas, estas podem mudar a qualquer momento do tempo. Isto pode ocasionar inconsistência em folksonomias aprendidas com dados de diferentes períodos de tempo. Assim, pode ser importante estudar o problema e



técnicas de solucionar o *concept drift* a fim de se evitar tais inconsistências.

# Referências Bibliográficas

- ABEL, F.; HENZE, N.; KRAUSE, D. A Novel Approach to Social Tagging: GroupMe!. In: 4<sup>th</sup> Int. Conf. on Web Information Systems and Technologies, pp. 42-49, 2008.
- ANDREWS, P.; PANE, J.; ZAIHRAYEU, I. Where are the Concepts in the Folksonomy Model? Technical Report # DISI-10-066, University of Trento. 2010.
- ASSAL, H.; SENG, J.; KURFESS, S.; SCHWARZ, E.; POHL, K. Partnering enhanced-NLP with semantic analysis in support of information extraction. In: Ontology-Driven Software Engineering, pp. 1-7, 2010.
- AUSTIN, J. L. How to do things with words. Oxford: Oxford University Press, p. 166. ISBN-13: 978-0674411524, 1962.
- BAO, S.; XUE, G.; WU, X.; YU, Y.; FEI, B.; SU, Z. Optimizing web search using social annotations. In: Proceedings of WWW '07, pp.501–510, 2007.
- BEAL, F.; WANDERLEY, G. M. P.; TACLA, C. A.; RAMOS, M. P.; PARAISO, E. C. FOLKUS-SD: Bulding Folksonomies From Source Code in Collaborative Software Development. In: Computer Supported Cooperative Work in Design, 2014.
- BEGELMAN, G., Keller, P., & SMADJA, F. Automated Tag Clustering: Improving search and exploration in the tag space. In: Collaborative Web Tagging Workshop at WWW'06, pp. 15-33, 2006.
- BOZ JR, G. ; RAMOS, M. P. ; SATO, G. ; NIEVOLA, J. C. ; PARAISO, E. C. . Noctua: a Tool for Knowledge Acquisition and Collaborative Knowledge Construction with a Virtual Catalyst. In: 15<sup>th</sup> International Conference on Computer Supported Cooperative

Work in Design (CSCWD). New York: IEEE. pp. 222-229, 2011.

BRIN, S.; PAGE, L. The Anatomy of a Large-Scale Hypertextual Web Search Engine. *Computer Networks and ISDN Systems*, 30(1-7), pp. 107–117, 1998.

CARLETTA, J.; ISARD, A.; ISARD, S.; KOWTKO, J.C.; DOHERTY-SNEDDON, G.; ANDERSON, A. H. The reliability of a dialogue structure coding scheme. *Computational Linguistics*, 23(1), pp. 13–31, 1997.

CATTUTO, C.; SCHMITZ, C.; BALDASSARRI, A.; SERVEDIO, V. D. P.; LORETO, V.; HOTHO, A.; GRAHL, M.; STUMME, G. Network properties of folksonomies. *AI Commun.*, 20(4), pp. 245–262, 2007.

CATTUTO, C.; BENZ, D.; HOTHO, A.; STUMME, G. Semantic Grounding of Tag Relatedness in Social Bookmarking Systems. In: *Proceedings of 7<sup>th</sup> International Conference on The Semantic Web*, pp. 1–16, 2008.

CHEN, W.; WANG, Y.; YANG, S. Efficient influence maximization in social networks. In: *Proc. of the 15<sup>th</sup> ACM Int. Conf. on Knowledge Discovery and Data Mining*, 2009.

CHOJNACKI, S.; KLOPOTEK, M. Random graph generative model for folksonomy network structure approximation. *Procedia Computer Science* 1(1), pp. 1683-1688, 2010.

DATTOLO, A.; PITASSI, E. Folkview: a multi-agent system approach to modeling folksonomies. In: *Proceedings of the 19<sup>th</sup> international conference on Advances in User Modeling*, pp. 198-212, 2012.

DAVIS, E. *Representations of Commonsense Knowledge*. Morgan Kaufmann Series in Representation and Reasoning. p. 515. Morgan Kaufmann Pub. ISBN-13: 978-1558600331. 1990.

- DI FELIPPO, A.; DIAS DA SILVA, B. C. Dos olhares sobre o léxico: diferenças e semelhanças. A construção de dicionários e bases de conhecimento lexical. Série Trilhas Linguísticas, n.9. Laboratório Editorial FCL/UNESP. Araraquara. Ed. Cultura Acadêmica, pp. 169-185. Disponível em: <[http://www.geterm.ufscar.br/ariani/Dos\\_olhares\\_sobre\\_o\\_lexico.pdf](http://www.geterm.ufscar.br/ariani/Dos_olhares_sobre_o_lexico.pdf)>. 2006.
- DICHEVA, D.; DICHEV, C. Can Collective Use Help for Searching?. International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery, pp. 24-31, 2011.
- ECHARTE, F.; ASTRAIN, J. J.; CÓRDOBA, A.; VILLADANGOS, J. Ontology of folksonomy: A New modeling method. In: Proceedings of Semantic Authoring, Annotation and Knowledge Markup, 2004.
- ECO, U. A theory of semiotics. Bloomington: Indiana University Press. ISBN: 0253359554 9780253359551. Disponível em: <<http://books.google.com.br/books?id=BVqwAAAAIAAJ>> p. 354, 1976.
- EMBLEY, D. W.; THALHEIM, B. Handbook of Concept Modeling: Theory, Practice, and Research Challenges. Springer, ISBN: 978-3-642-15865-0, p. 226, 2011.
- ERIKSSON, A. A survey of knowledge sources in dialogue systems. In: Proceedings of the IJCAI-99 Workshop on Knowledge and Reasoning in Practical Dialogue Systems. Murray Hill, New Jersey: International Joint Conference on Artificial Intelligence, pp. 41-48. Disponível em: <<http://citeseer.ist.psu.edu/328028.html>>. 1999.
- FARZANEH, H. H.; REIK, A. U.; MAURER, M. Development of a Knowledge Mapping Approach for Independent Knowledge Elicitation and Representation. In Impact: The Journal of Innovation Impact, inkt13-005: 5(1), pp. 43-53, 2013.
- FERREIRA, J. M. P. Folksonomias como conceitualizações compartilhadas na especificação

de modelos conceituais. Qualificação de doutorado, defesa em 16/10/2013. Programa de Pós-Graduação em Engenharia Elétrica e Informática Industrial. Universidade Tecnológica Federal do Paraná. Curitiba, 85 p., 2013.

GÓMEZ-PÉREZ, A.; FERNÁNDEZ-LÓPEZ, M.; CORCHO, O. *Ontological Engineering - Presentation*. Springer, (2004).

GRUBER, T. R. A translation approach to portable ontology specifications. *Knowledge Acquisition*, v.5, n.2, pp. 199-220. Academic Press. Disponível em: <<http://tomgruber.org/writing/ontolingua-kaj-1993.pdf>>. ISSN:1042-8143. 1993.

GRUBER, T. Collective knowledge systems: Where the Social Web meets the Semantic Web. *Web Semantics: Science, Services and Agents on the World Wide Web*, 06(1), pp. 4–13, 2007.

GUARINO, N.; OBERLE, D.; STAAB, S. *What is an Ontology?*. Dordrecht: Springer Verlag, 2009.

GUIZZARDI, G. *Ontological Foundations for Structural Conceptual Models*. University of Twente, Enschede. 2005.

GUPTA, M.; LI, R.; YIN, Z.; HAN, J. Survey on social tagging techniques. *SIGKDD Explor*, v.12, pp. 58–72, 2010.

HALPIN, H.; ROBU, V.; SHEPHERD, H. The complex dynamics of collaborative tagging. In: *Proceedings of the 16th International Conference on World Wide Web*, pp. 211–220, 2007.

HAZMAN, M.; EL-BELTAGY, S. R.; RAFAA, A. A survey of ontology learning approaches. *International Journal of Computer Applications*, 22(8), pp. 36–43, 2011.

HOTH, A.; JÄSCHKE, R.; SCHMITZ, C.; STUMME, G. *Information Retrieval in*

Folksonomies: Search and Ranking. In: *The Semantic Web: Research and Applications*, v. 4011 of LNAI, pp. 411–426, 2006a.

HOTHO, A.; JÄSCHKE, R.; SCHMITZ, C.; STUMME, G. Trend Detection in Folksonomies. In: *Prof. First International Conference on Semantics and Digital Media Technology*, v. 4306 of LNCS, pp. 56–70, 2006b.

HOTHO, A.; JÄSCHKE, R.; SCHMITZ, C.; STUMME, G. FolkRank: A Ranking Algorithm for Folksonomies. In: *Proc. FGIR*, pp. 111–114, 2006c.

IVANOVIC E. Automatic instant messaging dialogue using statistical models and dialogue acts. University of Melbourne. 2008.

JELASSI, M. N. A Quadratic Approach for Trend Detection in Folksonomies. *Web Reasoning and Rule Systems*. Vol. 7497, pp. 278-283, 2012.

JURAFSKY, D.; MARTIN, J. H. *Speech and Language Processing - An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. 2ed. p. 988. Upper Saddle River, Prentice Hall. Disponível em: <<http://www.cs.colorado.edu/~martin/SLP/Updates/1.pdf>>. ISBN-13: 978-0131873216. 2008.

KIM, H.; DECKER, S.; BRESLIN, J.; Representing and sharing folksonomies with semantics. *Journal of Information Science*. 36(1), pp. 57-72, 2010a.

KIM, H.; BRESLIN, J.; CHOI, J. H. Semantic representation for copyright metadata of user-generated content in folksonomies. *Online Information Review*. 34(4), pp. 626–641, 2010b.

KÖRNER, C.; BENZ, D.; HOTHO, A.; STROHMAIER, M.; STUMME, G. Stop Thinking, Start Tagging: Tag Semantics Emerge from Collaborative Verbosity. In: *Proceedings of*

the 19<sup>th</sup> International Conference on World Wide Web, pp. 521–530, 2010.

KNUTH, D. The Art of Computer Programming, Volume II: Seminumerical Algorithms. Addison-Wesley. 2<sup>nd</sup> Edition, ISBN: 0-201-03822-6, pp. 139-140, 1981.

LATORA, V.; MARCHIORI, M. Efficient behavior of small-world networks. Phys. Rev. Lett. Vol. 87, 198701. 2001.

LEE, K.; KIM, H.; SHIN, H.; KIM, H. J. Folksoviz: A semantic relation-based folksonomy visualization using the wikipedia corpus. Software Engineering, Artificial Intelligence, Networking, and Parallel/Distributed Computing, ACIS International Conference on, vol. 0, pp. 24-29, 2009.

LESTER, J.; BRANTING, K.; MOTT, B. Conversational agents. The Practical Handbook of Internet Computing. Chapman & Hall. ISBN-10: 9781584883814 8, PP. 2-3, 2004.

LIU, J.; GRUEN, D. M. Between Ontology and Folksonomy: A study of collaborative and implicit ontolgy evolution. In Proceedings of the 13th international conference on Intelligent user interfaces, pp. 361-364, 2008.

MATHES, A. Folksonomies - Cooperative Classification and Communication Through Shared Metadata. Library and Information Science, pp. 1–13, 2004.

MCTEAR, M. F. Spoken Dialogue Technology: Enabling the Conversational User Interface, New York. Journal CSUR ACM Computing Surveys. v. 34, issue. 1. Disponível em: <  
<http://www.ling.helsinki.fi/kit/2002s/ctl190net/materiaali/spokendialoguetechology.pdf>  
>. ISSN: 0360-0300. 2002.

MIKA, P. Ontologies are us: A unified model of social networks and semantics. Web Semantics: Science, Services and Agents on the World Wide Web, 5(1), pp. 5–15, 2007.

MILGRAM, S. The small world problem. *Psychol. Today* 2, pp. 60–67, 1967.

MINSKY, M. L. *A Framework for Representing Knowledge*. New York, MIT-AI Laboratory, Technical Report, Memo 306. Disponível em: <  
<http://web.media.mit.edu/~minsky/papers/Frames/frames.html> >. ISBN 55860-013-2.  
1975.

NEAL, J. G.; SHAPIRO, S. C. *Knowledge-Based Parsing*, Berlim, BOLC, L. (ed.) *Natural Language Parsing Systems*, pp. 49-92. Disponível em: <  
[http://link.springer.com/chapter/10.1007/978-3-642-83030-3\\_3#page-1](http://link.springer.com/chapter/10.1007/978-3-642-83030-3_3#page-1) >. ISBN: 978-3-642-83030-3. 1987.

PAULSEN, I.; MAINZ, D.; WELLER, K.; MAINZ, I.; KOHL, J.; VON HAESLER, A. *Ontoverse. Collaborative Knowledge Management in the Life Science Network*. In: *Proceedings of the German eScience Conference*, 2007.

PETERS, I. *Folksonomies: Indexing and retrieval in Web 2.0*. De Gruyter Saur, ISBN-10: 3598251793, 2009.

PLANGPRASOPCHOK, A.; LERMAN, K. *Constructing Folksonomies from User-Specified Relations on Flickr*. In: *Proceedings of the 18<sup>th</sup> International Conference on World Wide Web*, pp. 781-790, 2009.

PLANGPRASOPCHOK, A.; LERMAN, K.; GETOOR, L. *Growing a Tree in the Forest: Constructing Folksonomies by Integrating Structured Metadata*. In: *16<sup>th</sup> ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 949-958, 2010.

PLANGPRASOPCHOK, A.; LERMAN, K.; GETOOR, L. *A probabilistic approach for learning folksonomies from structured data*. In: *Fourth ACM International Conference on Web Search and Data Mining*, pp. 555–564, 2011.

QUINTARELLI, E. *Folksonomies: Power to the People*. In: *Proceedings of the ISKO Italy-*



UniMIB Meeting. Disponível em <<http://www.iskoi.org/doc/folksonomies.htm>>. 2005.

RUSSELL, S. J.; NORVIG, P. Artificial Intelligence: A Modern Approach (2 ed.). Pearson Education, ISBN: 85-352-1177-2, p. 769, 2003.

SCHMITZ, C.; HOTH, A.; JÄSCHKE, R.; STUMME, G. Mining Association Rules in Folksonomies. In: Lecture Notes in Computer Science - The Semantic Web: Research and Applications. Vol. 4011, pp. 411–426, 2006.

SEARLE, J. Speech Acts: An Essay in the Philosophy of Language. 1ed., p.203 Cambridge: Cambridge University Press. ISBN: 0-521-09626-X. Disponível em: <[http://books.google.com.br/books?hl=pt-BR&lr=&id=t3\\_WhfknvF0C&oi=fnd&pg=PA1&dq=+Speech+acts:+an+essay+in+the+philosophy+of+language&ots=0RmOaNR2Q-&sig=XUkYJhr9qu2O752vLQu2XbiS778#v=onepage&q=Speech%20acts%3A%20an%20essay%20in%20the%20philosophy%20of%20language&f=false](http://books.google.com.br/books?hl=pt-BR&lr=&id=t3_WhfknvF0C&oi=fnd&pg=PA1&dq=+Speech+acts:+an+essay+in+the+philosophy+of+language&ots=0RmOaNR2Q-&sig=XUkYJhr9qu2O752vLQu2XbiS778#v=onepage&q=Speech%20acts%3A%20an%20essay%20in%20the%20philosophy%20of%20language&f=false)>. 1969.

SEDGEWICK, R.; WAYNE, K. Algorithms (4 ed.). Addison-Wesley Professional, ISBN-10: 032157351X, pp.538-542, 2011.

SILVA, S. R.; BORTH, M. R.; FERREIRA, J. M. P.; FELTRIM, V. D. An approach to enrich users' personomy using the recommendation of semantic tags. Journal of the Brazilian Computer Society, v. 18, pp. 283-298, 2012.

SIPSER, M. Introdução à Teoria da Computação (2 ed.). Thomson Learning, ISBN: 978-85-221-0499-4, pp.10-12, 2007.

SPECIA, L.; MOTTA, E. Integrating Folksonomies with the Semantic Web. In: 4<sup>th</sup> European Semantic Web Conference. Vol. 4519, pp. 624–639, 2007.

SPITERI, L. The Structure and form of folksonomy tags: The road to the public library catalogue. Information Technology and Libraries, 27(3), pp. 13-25, 2007.

- STAAB, S.; SANTINI, S.; NACK, F.; STEELS, L.; MAEDCHE, A. Emergent semantics. *Intelligent Systems, IEEE*, 17(1), pp. 78–86, 2002.
- STEELS, L. The origins of ontologies and communication conventions in multi-agent systems. *Autonomous Agents and Multi-Agent Systems*, 1(2), pp. 169–194, 1998.
- STROHMAIER, M.; HELIC, D.; BENZ, D.; KÖRNER, C.; KERN, R. Evaluation of folksonomy induction algorithms. *ACM Trans. Intell. Syst. Technol.*, 2012.
- SZWARCFITER, J. L. *Grafos e Algoritmos Computacionais*. Editora Campus, ISBN: 85-7001-138-5, pp. 35-36, 1986.
- TAN, P.; STEINBACH, M.; KUMAR, V. *Introduction to Data Mining*. Addison-Wesley, ISBN-10: 0321321367, pp. 186-187, 2005.
- TACHIMORI, Y.; IWANAGA, H.; TAHARA, T. The networks from medical knowledge and clinical practice have small-world, scale-free, and hierarchical features. *Physica A: Statistical Mechanics and its Applications*, 392(23), pp. 6084-6089, 2013.
- TEMPICH, C.; PINTO, H. S.; SURE, Y.; STAAB, S. An argumentation Ontology for Distributed, Loosely-controlled and evolvinG Engineering processes of oNTologies (DILIGENT). In: *The Semantic Web: Research and Applications – Lecture Notes in Computer Science*, pp. 241–256, 2005.
- TRAUM, D. R.; HINKELMAN, E. A. Conversation acts in task-oriented spoken dialogue. *Computational Intelligence. Special Issue on Non-literal Language*, 8(3) 1992.
- TSYMBAL, A. The problem of concept drift: definitions and related work. Technical Report TCD CS-2004-15, Computer Science Department, Trinity College, Dublin. 2004.
- VAN DER WAL, T. Folksonomy Coinage and Definition. Disponível em:

<<http://vanderwal.net/folksonomy.html>> Acesso em: 02 fev. 2015. 2004.

- WANG, Y.; VOLKER, J.; HAASE, P. Towards semi-automatic ontology building supported by large-scale knowledge acquisition. In In AAI Fall Symposium on Semantic Web for Collaborative Knowledge Acquisition, pp.70-77, 2006.
- WATTS, D. J.; STROGATZ, S. H. Collective dynamics of 'small-world' networks. Nature 393, pp. 440–442, 1998.
- WEAVER, M. T.; FRANCE, R. K.; CHEN, Q.; FOX, E. A. A Frame-Based Language in Information Retrieval. Technical Report. Virginia Polytechnic Institute & State University, Blacksburg, VA, USA, 1988.
- WELLER, K. Folksonomies and Ontologies. Two New Players in Indexing and Knowledge Representation. In: H. Jezzard (Ed.), Applying Web 2.0. Innovation, Impact and Implementation. Online Information 2007 Conference Proceedings, pp. 108-115, 2007.
- WU C.; ZHOU, B. Tags are related: Measurement of semantic relatedness based on folksonomy network. Computing and Informatics, Vol. 30, pp. 165-188, 2011.
- WU, L.; YANG, L.; YU, N.; HUA, X. S. Learning to tag. In: Proceedings of the 18<sup>th</sup> international conference on World Wide Web, pp. 361– 370, 2009.
- XIAO, R.; NI, Y.; DU, X.; GONG, P. Towards a Correlation Cooccurrence Model Generating Approach to Folksonomy. In: Proceedings of the 2010 International Conference on Web Information Systems and Mining, Vol. 02, pp. 399-403, 2010.
- YOO, D.; SUH, Y. User-categorized tags to build a structured folksonomy. In: International Conference on Communication Software and Networks, pp. 160–164, 2010.
- ZUBIAGA, A. Enhancing Navigation on Wikipedia with Social Tags. In: Wikimania - 4<sup>th</sup> Annual Conference of the Wikimedia Community, 2009.