

ALDO MARCELO PAIM

**INFERÊNCIA DE PERSONALIDADE A PARTIR
DE TEXTOS EM PORTUGUÊS BRASILEIRO
UTILIZANDO LÉXICOS**

Dissertação 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

2016

ALDO MARCELO PAIM

**INFERÊNCIA DE PERSONALIDADE A PARTIR
DE TEXTOS EM PORTUGUÊS BRASILEIRO
UTILIZANDO LÉXICOS**

Dissertação 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. Fabrício Enembreck

CURITIBA

2016

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

P143 Paim, Aldo Marcelo
2016 Inferência de personalidade a partir de textos em português brasileiro
utilizando léxicos / Aldo Marcelo Paim; orientador, Fabrício Enembreck . -- 2016
160 f. : il. ; 30 cm

Dissertação (mestrado) – Pontifícia Universidade Católica do Paraná,
Curitiba, 2016
Bibliografia: f.123-142

1. Informática. 2. Algoritmos computacionais. 3. Mineração de dados
(Computação). 4. Personalidade. 5. Redes sociais. 6. Linguística. 7. Língua -
portuguesa – Lexicografia. I. Enembreck, Fabrício. II. Pontifícia Universidade
Católica do Paraná. Programa de Pós-Graduação
em Informática. III. Título.

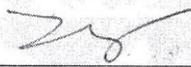
CDD 20. ed. – 004.068

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

DEFESA DE DISSERTAÇÃO DE MESTRADO Nº 18/2016

Aos 02 dias do mês de Março de 2016 realizou-se a sessão pública de Defesa da Dissertação "**Extração de Personalidade a partir de Textos em Português Brasileiro: utilizando Léxicos Linguísticos e Afetivos**" apresentado pelo aluno **Aldo Marcelo Paim**, como requisito parcial para a obtenção do título de Mestre em Informática, perante uma Banca Examinadora composta pelos seguintes membros:

Prof. Dr. Fabrício Enembreck
PUCPR (Orientador)


(assinatura)

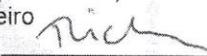
Aprovado
(Aprov/Reprov)

Prof. Dr. Emerson Cabrera Paraiso
PUCPR


(assinatura)

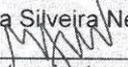
APROVADO
(Aprov/Reprov)

Prof. Dr. Richardson Ribeiro
UTFPR


(assinatura)

APROVADO
(Aprov/Reprov)

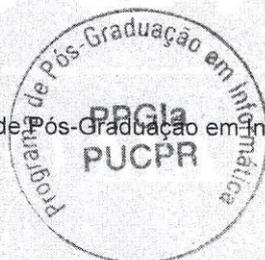
Prof.^a Dr.^a Maria Augusta Silveira Netto Nunes
UFS


(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 essa dissertação a meus pais, Alzenira e Aldemar.

Em especial a minha amada esposa Joana.

Agradecimentos

Agradeço primeiramente a Deus por ter me conduzido por esse caminho de grande amadurecimento, pela força e sabedoria para alcançar qualquer objetivo.

Em especial ao Prof. Dr. Fabrício Enembreck, aquele que permitiu meu ingresso na pesquisa. Sem sua ajuda, conselho, comentário e crítica, essa jornada não seria concluída. Aprendi muito com seus ensinamentos consequentes do vasto e admirável conhecimento. Obrigado pela confiança depositada e por transferir a mim qualidades inauditas que levarei à vida acadêmica.

Agradeço aos meus pais, que sempre me apoiaram e me proporcionaram condições para que eu pudesse chegar até aqui, me auxiliando de todas as formas e em todos os momentos.

A ela, minha amada esposa Joana, que esteve comigo nessa jornada sempre proferindo incentivos e não me deixando desanimar. Essa dissertação é dedicada a ti, juntamente com meus pais, pelas horas roubadas de seu convívio.

Quero agradecer também aos professores Emerson Paraíso, Maria Augusta Nunes e Richardson Ribeiro pelas contribuições e considerações feitas sobre esse trabalho. Suas colaborações foram de suma importância para a melhoria desse projeto como um todo.

Aos colegas de laboratório Jean, Heitor, André, Osmar, Ricardo e Jones por me auxiliarem nessa pesquisa com sugestões e revisões. Meu muito obrigado! Aprendi muito com vocês também.

Finalmente, quero agradecer à PUCPR, ao Programa de Pós-Graduação em Informática (PPGIa) pela oportunidade de realização de trabalhos em minha área de pesquisa e a todos que, de alguma forma, contribuíram para que esse trabalho fosse realizado.

Sumário

Agradecimentos	ii
Sumário	iii
Lista de Figuras	vi
Lista de Tabelas	viii
Lista de Símbolos	xiii
Lista de Abreviação	xv
Resumo	xvi
Abstract	xvii

Capítulo 1

Introdução	1
1.1. Motivação e Hipóteses.....	2
1.2. Objetivos.....	4
1.3. Organização.....	4

Capítulo 2

Fundamentação Teórica	5
2.1. Personalidade.....	5
2.1.1. Abordagem dos Traços.....	7
2.1.2. Modelo de 16 Fatores de Cattell.....	8
2.1.3. Modelo dos Três “Superfatores” de Eysenck.....	9
2.1.4. Modelo dos Cinco Grandes Fatores.....	10
2.1.5. Inventários de Personalidade do Modelo <i>BigFive</i>	12
2.2. Inferência de Personalidade por meio de Textos.....	15
2.2.1. Léxicos.....	21
2.2.2. Léxicos Afetivos.....	25
2.2.3. Outras abordagens	37
2.3. Inferência de Personalidade por meio de Textos para o Português do Brasil.....	41

2.4. Considerações Finais.....	46
Capítulo 3	
Um Modelo para Inferir a Personalidade a partir de Textos em Língua Portuguesa	48
3.1. Visão Geral do Método.....	48
3.2. Mensuração Explícita de Personalidade via Inventário.....	51
3.2.1. Participantes.....	53
3.3. Base de Dados.....	55
3.4. Pré-processamento.....	56
3.4.1. Utilização do Léxico LIWC.....	59
3.4.2. Utilização de TF-IDF	64
3.4.3. Utilização dos Léxicos Afetivos.....	66
3.5. Considerações Finais.....	75
Capítulo 4	
Procedimentos Metodológicos	77
4.1. Ferramentas de <i>Software</i> Utilizadas.....	77
4.1.1. Aplicação do Inventário.....	77
4.1.2. Coleta de Textos.....	78
4.1.3. Pré-Processamento.....	79
4.1.4. Waikato Environment for Knowledge Analysis – WEKA.....	80
4.2. Algoritmos de Aprendizagem de Máquina Utilizados.....	81
4.3. Avaliação dos Resultados.....	82
4.4. Considerações Finais.....	85
Capítulo 5	
Experimentos e Análise dos Resultados	86
5.1. Formação da Base de Dados Textual.....	86
5.2. Reconhecimento de Traços com TF-IDF.....	88
5.3. Reconhecimento de Traços com LIWC.....	93
5.4. Reconhecimento de Traços com LIWC Associado a Léxicos Afetivos.....	96
5.4.1. SentiStrength.....	96

5.4.2. AnewBr.....	101
5.4.3. Sentilex-PT.....	104
5.4.4. OpLexicon.....	108
5.4.5. Experimentos com a União dos Léxicos Afetivos	110
5.5. Reconhecimento de Traços com LIWC Associado a Léxicos Afetivos e TF-IDF.....	113
5.6. Análise dos Resultados.....	115
5.7. Considerações Finais.....	118
Capítulo 6	
Conclusão e Trabalhos Futuros	120
Referências Bibliográficas	123
Apêndice A	
Termo de consentimento para uso de informações pessoais	143
Apêndice B	
Resultados Suplementares de Experimentos	146
Anexo 1	
Questões do Inventário NEO-IPIP 120	152
Anexo 2	
Parecer consubstanciado do Comitê de Ética em Pesquisa da PUCPR	157

Lista de Figuras

Figura 2.1	Divisão da abordagem baseado em léxico para a Análise de Sentimento em textos. Adaptada de (MEDHAT; HASSAN; KORASHY, 2014).	26
Figura 2.2	Exemplo de classificação SentiWordNet. Adaptada de (SENTIWORDNET, 2015).....	28
Figura 2.3	Escala de avaliação de valência (A) e alerta (B) do <i>Self-Assessment Manikin</i> . (KRISTENSEN <i>et al.</i> 2011).....	32
Figura 2.4	Amostragem do arquivo CSV do SentiLex-PT (SENTILEX-PT, 2015).....	35
Figura 2.5	Processo de classificação semi-supervisionado utilizado por (LIMA; CASTRO, 2013).....	45
Figura 3.1	Visão geral do modelo para inferir a personalidade a partir de textos....	50
Figura 3.2	Questão 1 do Teste NEO-IPIP 120. Adaptado de (JOHNSON, 2014)...	52
Figura 3.3	Média dos resultados com desvio padrão dos traços de personalidade via inventário.....	55
Figura 3.4	Etapas do pré-processamento.....	58
Figura 3.5	LIWC na etapa do pré-processamento.....	60
Figura 3.6	Demonstração do vetor de características.....	62
Figura 3.7	Estrutura da base de treinamento.....	63
Figura 3.8	Estrutura da base de treinamento com os atributos LIWC e TF-IDF.....	65
Figura 3.9	Estrutura da base de treinamento com a contabilização das emoções do léxico <i>SentiStrength</i>	67
Figura 3.10	Estrutura da base de treinamento com a contabilização das emoções do léxico <i>SentiStrength</i> baseado na frequência dos termos.....	69
Figura 3.11	Estrutura da base de treinamento com a contabilização das emoções do léxico <i>AnewBr</i> baseado na frequência do termo.....	71
Figura 3.12	Estrutura da base de treinamento com a contabilização do <i>SentiLex-PT</i> em duas abordagens: (i) soma dos termos emotivos; (ii) baseada na frequência dos termos.....	72

Figura 3.13	Estrutura da base de treinamento com a contabilização do <i>OpLexicon</i> em duas abordagens: (i) soma dos termos emotivos; (ii) baseada na frequência dos termos.....	74
Figura 3.14	Demonstração da base de treinamento com a combinação dos léxicos: LIWC, <i>OpLexicon</i> e <i>AnewBr</i>	75
Figura 4.1	Tela principal do software responsável pela etapa de pré-processamento dos dados.....	80
Figura 5.1	Resultado do teste de Nemenyi para o comparativo dos algoritmos do experimento TF-IDF para os traços: Neuroticismo e Socialização.....	92
Figura 5.2	Resultado do teste de Nemenyi para o comparativo dos algoritmos do experimento LIWC associado ao léxico <i>SentiStrength</i>	100
Figura 5.3	Resultado do teste de Nemenyi para o comparativo dos algoritmos do experimento LIWC associado ao léxico <i>AnewBr</i>	104
Figura 5.4	Resultado do teste de Nemenyi para o comparativo dos algoritmos do experimento LIWC associado ao léxico <i>OpLexicon</i>	110
Figura 5.5	Resultado do teste de Nemenyi para o comparativo dos algoritmos do experimento LIWC associado a todos os léxicos afetivos.....	112

Lista de Tabelas

Tabela 2.1	Os 16 fatores proposto por Cattell (SOUZA; PRIMI; MIGUEL, 2007) ..	8
Tabela 2.2	Adjetivos característicos dos Cinco Grandes Fatores (SOUZA; NUNES, 2011).....	11
Tabela 2.3	Facetas do questionário NEO-PI-R (NETO, 2010).....	12
Tabela 2.4	Quantidade de questões dos inventários baseados no modelo <i>BigFive</i>	13
Tabela 2.5	Questões do inventário NEO-IPIP 300 baseados no modelo <i>BigFive</i> , adaptado de (NUNES, 2008).....	14
Tabela 2.6	Comparação dos melhores modelos para cada traço. (1) Precisão da classificação; (2) porcentagem de melhoria em relação à <i>baseline</i> utilizando regressão; (3) perda de ranking. Adaptado de (MAIRESSE et al., 2007).....	18
Tabela 2.7	Características textuais do LIWC, adaptado de (PENNEBAKER et al., 2007).....	22
Tabela 2.8	Dimensões e categorias do LIWC para o português (BALAGE FILHO et al., 2013).....	24
Tabela 2.9	Lista de “ <i>a-labels</i> ” com os respectivos estados afetivos e exemplos (PASQUALOTTI; VIEIRA, 2008).....	29
Tabela 2.10	Estrutura da <i>WordnetAffect</i> e exemplos de registros da base (PASQUALOTTI; VIEIRA, 2008).....	30
Tabela 2.11	Amostra do léxico <i>EmoSenticNet</i> (EMOSENTICNET, 2015).....	33
Tabela 2.12	Estrutura do léxico <i>OpLexicon</i> (OPLEXICON, 2015).....	36
Tabela 2.13	Resumo do estado da arte das abordagens de identificação de personalidade por meio de texto.....	39
Tabela 2.14	Exemplo de termos do dicionário produzido pelos autores (NUNES; TELES; DE SOUZA, 2013).....	42
Tabela 2.15	Correlação entre NEO-IPIP & Text-Mining, adaptado de (NUNES; TELES; DE SOUZA, 2013).....	43
Tabela 2.16	Correlação entre TIPI & Text-Mining, adaptado de (NUNES; TELES; DE SOUZA, 2013).....	43

Tabela 2.17	Pré-processamento dos grupos, formando o conjunto meta-base com 30 objetos e 11 atributos, adaptado de (LIMA; CASTRO, 2013).....	46
Tabela 2.18	Amostra de classificação do ppBates, adaptado de (LIMA; CASTRO, 2013)	46
Tabela 3.1	Informações sobre os participantes do experimento.....	53
Tabela 3.2	Informações sobre os participantes finais do experimento.....	52
Tabela 3.3	Exemplos de <i>stopwords</i>	57
Tabela 3.4	Categorias do LIWC para o português utilizadas no experimento.....	62
Tabela 3.5	Representação das publicações usando peso TF-IDF.....	64
Tabela 4.1	Valores de referência para a interpretação do coeficiente de correlação Pearson. Adaptado de (APPOLINÁRIO, 2006).....	84
Tabela 5.1	Amostra dos textos coletados na rede social <i>Facebook</i>	87
Tabela 5.2	Informações sobre o conjunto de dados.....	87
Tabela 5.3	Conjunto de termos TF-IDF extraídos dos textos.....	88
Tabela 5.4	Lista de 68 termos extraídos por meio de TF-IDF nas publicações dos usuários.....	89
Tabela 5.5	Resultados de correlação de Pearson do experimento TF-IDF para o traço Extroversão (* correlação significativa ao nível de 0,05).....	89
Tabela 5.6	Resultados de correlação de Pearson do experimento TF-IDF para o traço Neuroticismo (* correlação significativa ao nível de 0,05).....	90
Tabela 5.7	Resultados de correlação de Pearson do experimento TF-IDF para o traço Realização (* correlação significativa ao nível de 0,05).....	90
Tabela 5.8	Resultados de correlação de Pearson do experimento TF-IDF para o traço Socialização (* correlação significativa ao nível de 0,05).....	90
Tabela 5.9	Resultados de correlação de Pearson do experimento TF-IDF para o traço Abertura (* correlação significativa ao nível de 0,05).....	91
Tabela 5.10	Melhores correlações TF-IDF com os traços de personalidade.....	91
Tabela 5.11	Resultados de correlação de Pearson do experimento LIWC para todos os traços (* correlação significativa ao nível de 0,05).....	94
Tabela 5.12	Resultados de correlação de Pearson do experimento LIWC com seletor de atributos para todos os traços (* correlação significativa ao nível de	

	0,05).....	95
Tabela 5.13	RMSE do experimento LIWC com seletor de atributos para todos os traços.....	95
Tabela 5.14	Estatística das polaridades do léxico <i>SentiStrength</i> identificados nos textos.....	97
Tabela 5.15	Resultados de correlação de Pearson do experimento da primeira abordagem dos léxicos LIWC e <i>SentiStrength</i> (* correlação significativa ao nível de 0,05).....	97
Tabela 5.16	Comparativo entre os resultados de correlação de Pearson entre LIWC e LIWC com <i>SentiStrength</i> , utilizando sumarização de polaridades (* correlação significativa ao nível de 0,05).....	98
Tabela 5.17	Lista de termos TF-IDF identificados no léxico <i>SentiStrength</i>	99
Tabela 5.18	Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e <i>SentiStrength</i> (* correlação significativa ao nível de 0,05).....	99
Tabela 5.19	Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e <i>SentiStrength</i> com seleção de atributos (* correlação significativa ao nível de 0,05).....	100
Tabela 5.20	Resultados de correlação de Pearson do experimento da primeira abordagem dos léxicos LIWC e <i>AnewBr</i> (* correlação significativa ao nível de 0,05).....	101
Tabela 5.21	Lista de termos TF-IDF identificados no léxico <i>AnewBr</i>	102
Tabela 5.22	Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e <i>AnewBr</i> (* correlação significativa ao nível de 0,05).....	103
Tabela 5.23	Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e <i>AnewBr</i> com seleção de atributos (* correlação significativa ao nível de 0,05).....	103
Tabela 5.24	Estatística das polaridades do léxico <i>SentiLex-PT</i> identificados nos textos.....	105
Tabela 5.25	Resultados de correlação de Pearson do experimento da primeira abordagem dos léxicos LIWC e <i>SentiLex-PT</i> (* correlação significativa	105

	ao nível de 0,05).....	
Tabela 5.26	Lista de termos TF-IDF identificados no léxico <i>SentiLex-PT</i>	106
Tabela 5.27	Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e <i>SentiLex-PT</i> (* correlação significativa ao nível de 0,05).....	107
Tabela 5.28	Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e <i>SentiLex-PT</i> com seleção de atributos. (* correlação significativa ao nível de 0,05).....	107
Tabela 5.29	Resultados de correlação de Pearson do experimento da primeira abordagem dos léxicos LIWC e <i>OpLexicon</i> (* correlação significativa ao nível de 0,05).....	108
Tabela 5.30	Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e <i>OpLexicon</i> (* correlação significativa ao nível de 0,05).....	109
Tabela 5.31	Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e <i>OpLexicon</i> com seleção de atributos (* correlação significativa ao nível de 0,05).....	110
Tabela 5.32	Correlação de Pearson do experimento com todos os léxicos (* correlação significativa ao nível de 0,05).....	111
Tabela 5.33	Correlação de Pearson do experimento de seleção de atributos da base de dados com todos os léxicos (* correlação significativa ao nível de 0,05).....	112
Tabela 5.34	Correlação de Pearson do experimento combinando LIWC, léxicos afetivos e TF-IDF-750, sem seleção de atributos (* correlação significativa ao nível de 0,05).....	113
Tabela 5.35	Correlação de Pearson e RSME do experimento combinando LIWC, léxicos afetivos e TF-IDF-750, com seleção de atributos (* correlação significativa ao nível de 0,05).....	114
Tabela 5.36	Correlação de Pearson do experimento combinando LIWC, léxicos afetivos e TF-IDF-1000, sem seleção de atributos (* correlação	

	significante ao nível de 0,05).....	114
Tabela 5.37	Correlação de Pearson e RSME do experimento combinando LIWC, léxicos afetivos e TF-IDF-1000, com seleção de atributos (* correlação significante ao nível de 0,05).....	115
Tabela 5.38	Comparação dos resultados de correlação de Pearson entre os léxicos afetivos que utilizaram a sumarização de polaridades.....	117
Tabela 5.39	Comparação dos resultados de correlação de Pearson entre os léxicos afetivos que utilizaram peso de frequência (TF-IDF).....	117

Lista de Símbolos

w	Total de palavras já utilizadas em uma sentença
T	Total de palavras em uma sentença
x_i	Quantidade total de palavras do usuário
s_1	Fator de corte com exíguas quantidades de palavras
s_2	Fator de corte com excessivas quantidades de palavras
#	<i>Hashtags</i>
M_v	Média ponderada para a valência
M_A	Média ponderada para a alerta
q_i	Quantidade de vezes que uma palavra i é encontrada
v_i	Valor de valência de uma palavra i
A_i	Valor de alerta de uma palavra i
t	Termo
d	Documento
l	Léxico
$TFIDF(t, d)$	Peso de frequência do termo t em um documento d
$P_{t,l}$	Polaridade do termo t no léxico l .
r	Coefficiente de correlação
n	Número de elementos no vetor
x	Vetor de valores dos dados reais
y	Vetor dos valores dos dados obtidos
i	Representa o i -ésimo elemento do vetor
\bar{x}	Média dos valores do vetor x
\bar{y}	Média dos vetores de y
α	Nível de significância para testes estatísticos
$LIWC_{(t,d)}$	Valor do léxico LIWC para o termo t em um documento d
$tf_{t,d}$	Ocorrência do termo t em um documento d

tt_d

Total de palavras empregadas no documento d

Lista de Abreviações

AM	Aprendizagem de Máquina
AS	Análise de Sentimentos
RPT	Reconhecimento de Personalidade por meio de Texto
BI	<i>Business Intelligence</i>
TIPI	<i>Ten-Item Personality Inventory</i>
FFPI	<i>Five Factor Personality Inventory</i>
NEO-IPIP	<i>Neo-International Personality Item Pool</i>
BFQ	<i>Big Five Questionnaire</i>
NEO-PI-R	<i>Revised NEO Personality Inventory</i>
GPI	<i>Global Personality Inventory</i>
LIWC	<i>Linguistic Inquiry and Word Count</i>
MRC	<i>Medical Research Council</i>
ANEW	<i>Affective Norms for English Words</i>
ANEWBR	<i>Brazilian norms for the Affective Norms for English Words</i>
SAM	<i>Self-Assessment Manikin</i>
API	<i>Application Programming Interface</i>
HTTP	<i>Hypertext Transfer Protocol</i>
JSON	<i>JavaScript Object Notation</i>
KNN	<i>K-Nearest Neighbors</i>
TF-IDF	<i>Term Frequency - Inverse Document Frequency</i>
RMSE	<i>Root Mean Squared Error</i>
NB	<i>Naïve Bayes</i>
LR	<i>Linear Regression</i>
SVM	<i>Support Vector Machine</i>

Resumo

Os avanços recentes em análise automática de textos fomentaram o surgimento de uma área de pesquisa responsável por reconhecimento de aspectos subjetivos, tais como personalidade, opiniões e emoções que os autores empregam em seus textos. Pesquisas nessa área remetem ao desenvolvimento de métodos que possibilitam que sistemas computacionais sejam capazes de reconhecer e detectar características de personalidade em texto. Entretanto, por ser uma área relativamente nova, esses métodos ainda estão em fase de desenvolvimento e são, em sua grande maioria, para a língua inglesa. Desta forma, observa-se a necessidade de desenvolvimento de técnicas com o mesmo objetivo para outros idiomas, como por exemplo, a língua portuguesa. Neste estudo é apresentado um método para a inferência da personalidade do indivíduo por meio de textos publicados em redes sociais, escritos em língua portuguesa. O método é baseado em características linguísticas e afetivas, por intermédio de léxicos. Essas informações são processadas por meio de algoritmos de mineração de textos para geração de modelos que descrevem os cinco grandes fatores da personalidade de um indivíduo (*BigFive*), segundo a teoria de traços: Extroversão, Neuroticismo, Realização, Socialização e Abertura à experiência. Os resultados apresentam valores de saída fortemente correlacionados à personalidade humana para todos os traços.

Palavras-Chave: Mineração de Texto; Computação Afetiva; Reconhecimento de Personalidade; Análise de Redes Sociais.

Abstract

Recent advances in automatic text analysis fostered the emergence of an area responsible for recognition of subjective aspects such as personality, opinions and emotions that the authors use in its texts. Research in this area refer to the development of methods that allow computer systems to recognize and detect personality features into text. However, as it is a relatively new area, these methods are still in the development phase and are, in its vast majority, only for the English language. Thus, there is a development need of techniques aiming other languages, such as Portuguese. In this study we present a proposal to create a method for the inference of the individual's personality through texts published on social networks, written in Portuguese. The method is based on linguistic features and affectivity, both through lexicons. This information is processed through the text mining algorithms for generating models that describe the five major factors of an individual's personality (*BigFive*), according to the theory of traits: Extraversion, Neuroticism, Conscientiousness, Agreeableness and Openness to Experience. Results show output values strongly correlated with human personality for all traits.

Keywords: Text Mining; Affective Computing; Personality Recognition; Social Networks Analysis.

Capítulo 1

Introdução

A descoberta do comportamento e esclarecimento da personalidade humana está presente em diversas áreas do conhecimento: Psicologia, Filosofia, Sociologia, Antropologia e Medicina em geral. A tendência em classificar pessoas é tão antiga quanto a humanidade, uma vez que muitas características da personalidade dos indivíduos estão relacionadas ao meio social e cultural em que vivem (MATHEWS et al., 2009).

Nos últimos anos a inferência de personalidade humana tem sido objeto de pesquisa nas áreas da computação. A *World Wide Web* teve um fator fundamental para tais estudos, pois no mundo virtual cada vez mais relações sociais acontecem, o que possibilita que os usuários expressem suas opiniões e discutam suas ideias, demonstrando até mesmo posicionamentos pessoais diante de determinadas situações. A natureza afetiva desses pareceres gradativamente se torna a base para a tomada de decisões sobre pesquisa de *marketing*, BI (*Business Intelligence*), previsão no mercado de ações e monitoramento de imagem (MONTROYO et al., 2012).

Dessa maneira, a computação permitiu a análise de grandes quantidades de texto com o intuito de descobrir de forma automática os traços de personalidade de seus autores, denominando-se tal processo como Reconhecimento da Personalidade através de Texto (RPT doravante), do inglês *Personality Recognition from Text* (CELLI, 2012). Esta tarefa, que é parcialmente ligada à atribuição da autoria, requer habilidades e técnicas de várias áreas diferentes, como Linguística, Psicologia e Ciências da Comunicação (CELLI, 2012). Diversos pesquisadores, como (MAIRESSE, 2007; QUERCIA et al., 2012; BACHRACH et al., 2012; CELLI et al., 2014), utilizam o modelo de fatores chamado “*BigFive*”, que descreve a personalidade de um indivíduo em torno de cinco traços (NORMAN, 1963):

- Extroversão (sociável, ativo e assertivo)
- Neuroticismo (neurótico)
- Socialização (amigável e cooperativos)
- Realização (organizado e disciplinado)
- Abertura à experiência (intelectual e aberto ao novo).

Em relação ao RPT, a comunidade científica tem se concentrado principalmente em expandir sua aplicação para idiomas diferentes do inglês, como por exemplo, os estudos de (BAI et al., 2012) e (KERMANIDIS, 2012). No caso da língua portuguesa essa expansão é limitada, existindo escassos estudos computacionais com o objetivo de inferir a personalidade humana (NUNES; TELES; DE SOUZA, 2013; LIMA; CASTRO, 2013).

Devido à carência de trabalhos que tratam da identificação de personalidade, a presente pesquisa visa promover um modelo de reconhecimento automático de personalidade a partir de texto em língua portuguesa, através de uma abordagem léxica associada a um método de representação de termos com base em sua frequência (*Term Frequency - Inverse Document Frequency*), utilizando algoritmos de aprendizagem de máquina (AM). Ainda, o estudo conta com a colaboração de psicólogo para o desenvolvimento do método proposto.

1.1. Motivação e Hipóteses

Estudos desenvolvidos no campo da Neurociência e da Psicologia (DAMASIO, 1996; TRAPPL et al., 2003; THAGARD, 2006) comprovam o papel fundamental que a personalidade do indivíduo reflete na cognição humana, no que se refere à percepção, raciocínio, aprendizagem, memória e tomada de decisão. Aspectos sutis na interação homem-máquina podem revelar facetas da personalidade humana, já que muitas vezes os indivíduos respondem psicologicamente a computadores pensando que eles também são humanos (REEVES; NASS, 1996).

Dessa maneira, a *Internet* tem desempenhado um papel fundamental na referida interação, se tornando um ambiente com enormes repositórios de dados escritos, adequado para o reconhecimento da personalidade. Usuários presentes em redes sociais, *blogs* e demais meios que permitem expressar opinião, gerar informações e discutir ideias, tornam-se dados relevantes para a indústria. A inferência de personalidade em tais meios de comunicação melhora expressivamente a personalização de produtos, otimização de serviços e a tomada de decisão.

Diante da importância dessa aplicação à área da computação, o presente estudo tem o intuito de propor um modelo computacional capaz de revelar traços de personalidade por meio de textos em português (brasileiro) publicados em mídias sociais, utilizando léxicos psicolinguísticos, afetivos e um método de representação de termos com base em sua frequência. A utilização de ambos os léxicos torna-se um diferencial para essa pesquisa, pois explora a relação entre a personalidade e o aspecto emocional.

Observa-se que os métodos de inferência de personalidade para o português brasileiro existentes na literatura, utilizam limitadas características textuais (LIMA; CASTRO, 2013), abstém as emoções empregadas nos textos e avaliam o método a uma exígua quantidade de participantes (NUNES; TELES; DE SOUZA, 2013). Essas limitações podem ser parcialmente explicadas, pois se obter um modelo de inferência de personalidade não é trivial, especialmente em ambientes virtuais, onde ocasionalmente publicações estão protegidas por regras de privacidade, além de requerer conhecimentos sólidos em Linguística, Psicologia e Mineração de Dados.

O modelo proposto utilizando uma abordagem de léxicos é motivado por dois aspectos importantes: (i) um léxico com elevado número de características a serem extraídas de um texto promove maior desempenho no reconhecimento da personalidade; e (ii) a identificação de emoções contribui com a averiguação dos traços de personalidade, uma vez que a personalidade é identificada por características de um organismo autônomo que representa padrões consistentemente escolhidos de reação mental, incluindo emoções (MOFFAT, 1997).

Outro fator motivacional determinante à pesquisa é a expansão do RPT para a língua portuguesa, considerando que existe carência de trabalhos que tratam da identificação de personalidade por meio de textos na área da computação para essa língua, tendo em vista que a grande maioria dos estudos se concentra no idioma inglês.

As hipóteses deste trabalho são: (i) que o reconhecimento da personalidade pode ser mensurado através de trechos escritos por seus autores, (ii) devido à lacuna de métodos para a inferência em língua portuguesa, uma abordagem linguística com muitas características textuais pode agregar maior desempenho no reconhecimento da personalidade, (iii) mostrar que através do acréscimo de um léxico afetivo e uma abordagem de representação de termos com base em sua frequência, é possível tornar o modelo robusto e com melhores resultados.

1.2. Objetivos

Este trabalho tem como objetivo principal o desenvolvimento de um método para inferência da personalidade humana a partir de textos publicados *online*, escritos em língua portuguesa, por meio da análise de texto. Os objetivos específicos incluem o levantamento bibliográfico do reconhecimento de personalidade a partir de textos, o estudo de métricas de *text-mining* utilizando léxicos disponíveis em português, a mensuração dos traços de personalidade dos participantes aplicando inventários explícitos, a coleta dos textos públicos dos participantes e implementação de um modelo que possibilite a mineração desses textos inferindo traços de personalidade por meio de pistas deixadas pelos autores, assim como uma avaliação empírica dos resultados, utilizando testes de hipótese não paramétricos.

1.3. Organização

Este trabalho está organizado da seguinte maneira: no Capítulo 2 descreve-se os principais conceitos sobre personalidade humana, apresenta-se estudos na área da computação afetiva, conceitua o tema e suas abordagens, posteriormente se discute os avanços em diversas línguas, inclusive a língua portuguesa. No Capítulo 3 apresenta-se as etapas de construção do método proposto para a inferência de personalidade por meio de textos em português. Em seguida, no Capítulo 4 cita-se as ferramentas e as tecnologias usadas na implementação do método descrito no Capítulo 3. No Capítulo 5 apresenta-se os principais experimentos realizados com o método de identificação de personalidade e faz uma análise dos resultados obtidos. Por fim, serão apresentadas, no Capítulo 6, as conclusões e propostas para trabalhos futuros.

Capítulo 2

Fundamentação Teórica

No presente Capítulo são apresentados os conceitos que estão relacionados com o trabalho de reconhecimento de personalidade a partir de texto. Preliminarmente, serão mostrados alguns pressupostos conceituais de suma importância para compreender o assunto principal da pesquisa. Neste contexto, as primeiras seções do Capítulo apresentam a definição de personalidade, bem como a abordagem dos traços de personalidade, expondo os principais modelos. Em seguida, serão analisados os conceitos de reconhecimento de personalidade por meio de texto, assim como os principais estudos da área, abordando a utilização de léxicos psicolinguísticos, afetivos e outros métodos. Por fim, a seção seguinte discute o estudo de reconhecimento de personalidade a partir de texto para o português do Brasil.

2.1. Personalidade

A palavra personalidade deriva do latim - *persona* - que significa máscara, isto é, como nos apresentamos aos outros indivíduos da sociedade (BAISE, 2008). A Psicologia possui dessemelhantes definições para personalidade humana. Em (TAVARES, 2006), o autor define a personalidade como um conjunto de características que determina como o indivíduo pensa, sente e age, conferindo-lhe uma identidade e um padrão de relacionamentos único. Essas características são forjadas pela interação entre disposições naturais e experiências ao longo do tempo, permitindo flexibilidade e um melhor ajuste do indivíduo ao ambiente. Schultz (1990) define a personalidade como um conjunto de aspectos internos e externos peculiares relativamente permanentes do caráter de uma pessoa que influenciam o comportamento em diferentes situações.

Observa-se que uma das definições de personalidade mais mencionadas em textos didáticos e científicos da Psicologia é a de Gordon Allport (ALLPORT 1961), como por exemplo, citado nas obras de (BAISE, 2008; HENNA, 2011; URSO, 2011; DWEEK, 2008; ALMEIDA, 2010). Allport afirma no livro *Personality: a psychological interpretation*, do ponto de vista da Psicologia, que a personalidade refere-se ao "que o homem realmente é", e complementa: "a organização interna e dinâmica de sistemas psicóticos do indivíduo que determinam o seu ajuste único ao ambiente". Allport foi um dos pioneiros no campo da Psicologia no aprendizado da personalidade do indivíduo.

Apesar de inúmeras citações e de ter um importante papel histórico, há autores que declaram que a definição de Allport não é correlacionada com o entendimento contemporâneo da personalidade (PERVIN; JOHN, 2003).

De fato, a definição da personalidade é demasiadamente complexa, talvez a definição compreendida pela maioria dos autores e teóricos da personalidade é que esta pode ser associada a um sistema, cujo conjunto de padrões inatos de cada indivíduo interage com o ambiente social nas dimensões afetivas, cognitivas e comportamentais para produzir as ações e as experiências de uma vida individual, variando e dependendo da abordagem (ou teoria) utilizada (GARCIA, 2006).

As principais abordagens de personalidades humanas são listadas a seguir (NETO, 2010):

- Abordagem psicanalítica: Tal abordagem examina como as forças do inconsciente, desejos, motivos e necessidades físicas e biológicas induzem os sentimentos, pensamentos e comportamentos.
- Abordagem humanista: Essa abordagem foca na experiência pessoal, sentimentos e valores. Tal proposta formula a personalidade de cada indivíduo através da interpretação e observação do mundo. Psicólogos que utilizam essa abordagem acreditam que a experiência do indivíduo é subjetiva e que as pessoas são genuinamente boas.
- Abordagem dos traços: Baseia-se nas palavras utilizadas para descrever a personalidade do indivíduo. Esta abordagem procura identificar quais traços descrevem melhor uma pessoa e quanto eles diferem para outrem.

- Abordagem cognitiva: Enfatiza como o comportamento e a personalidade são moldados com as crenças, experiências prévias, histórias individuais e interações com o ambiente.
- Abordagem Biológica: Apresenta os fatores genéticos e neuroquímicos que explicam a predisposição de indivíduos a determinados comportamentos.

Das definições acima descritas, a abordagem de traços foi a que mais influenciou o desenvolvimento de testes e modelos de personalidade (NUNES, 2012). Isso é explicado pelo nível de aprofundamento descritivo que tal abordagem se propõe a fazer, e pela semelhança com que as pessoas percebem de maneira intuitiva a personalidade (NETO, 2010). A seção a seguir descreve com maiores detalhes a abordagem dos traços e modelos de personalidade.

2.1.1. Abordagem dos Traços

A abordagem dos traços assemelha-se da maneira intuitiva a como uma pessoa descreve a personalidade de outro indivíduo, isto é, para uma pessoa descrever a personalidade de outra, comumente não se detém a conhecimentos biológicos, ou ainda aprofundamento em conflitos de inconsciente, mas provavelmente descreverá uma pessoa como sendo otimista ou pessimista, introvertida ou extrovertida ou qualquer outra característica de personalidade. No contexto da abordagem dos traços, cada uma dessas características é chamada de traço. Segundo Urso (URSO, 2006) um traço é uma característica individual temporalmente estável, expresso por uma coleção de comportamentos, atitudes e emoções repetitivos e habituais mesmo em diferentes circunstâncias e situações interpessoais.

Os pesquisadores da universidade de Harvard, Gordon Allport e Henry Obdert foram os primeiros pesquisadores a trabalhar o conceito de traços psicológicos para descrever a personalidade humana, desenvolvendo a teoria dos traços (ALLPORT; ALLPORT, 1921). Para os autores, cada pessoa possui traços de personalidade comuns e individuais, e ainda traços latentes que permitem descrever a personalidade de um indivíduo.

Allport e Obdert percorreram o dicionário e assinalaram todas as palavras que poderiam ser utilizadas para descrever traços de personalidade. Como resultado, 17.953 palavras foram selecionadas – traços comuns e individuais – na sua maioria adjetivos. Esse vocábulo foi reduzido a pouco mais de 4.500 após remoção de sinônimos (NETO, 2010).

Vários pesquisadores se inspiraram no trabalho de Allport e Odbert e criaram diversos modelos, para classificar a personalidade de um indivíduo utilizando um determinado número de traços. Os principais modelos desenvolvidos apresentados nesse trabalho são: (i) “Modelo dos 16 Fatores”, criado por Raymond Cattell para avaliar a personalidade de um indivíduo em relação aos 16 fatores primários; (ii) “Modelo dos Três Superfatores”, de Hans Eysenck e; (iii) “Modelo dos Cinco Fatores”. Os detalhes de cada modelo serão apresentados a seguir.

2.1.2. Modelo de 16 Fatores de Cattell

Raymond Bernard Cattell, psicólogo britânico, expandiu o desenvolvimento dos traços de personalidade criados por Allport e Odbert e criou a teoria de traços utilizando a análise fatorial, em uma tentativa de criar uma “tabela periódica” dos elementos da personalidade (NUNES, 2012).

A análise fatorial é um procedimento estatístico baseado no conceito de coeficiente de correlação que mensura o relacionamento entre dois conjuntos de variáveis (NETO, 2010). Em 1946, Cattell foi um dos pioneiros a utilizar a computação para percorrer a extensa base de Allport, a fim de agrupar os descritores de forma objetiva. Como resultado final, Cattell construiu um modelo de dezesseis fatores primários, propondo a hipótese que esses fatores independentes são capazes de descrever as principais características da personalidade humana. Os fatores refinados por Cattell são apresentados na tabela 2.1.

Tabela 2.1: Os 16 fatores proposto por Cattell (SOUZA; PRIMI; MIGUEL, 2007).

Fatores	Valores Baixos	Valores Altos
Expansividade	reservado, impessoal, distante	expansivo, participante, atencioso
Inteligência	menos inteligente, pensamento concreto	mais inteligente, pensamento abstrato
Estabilidade	emocional, sensível às impressões afetivas, emocionalmente instável	emocionalmente estável, adaptável, maduro
Afirmação	hulmidade, brando, cooperativo, avesso a conflitos	afirmativo, dominante, agressivo, assertivo
Preocupação	sóbrio, sério, retraído, prudente	despreocupado, alegre, animado
Consciência	evasivo, inconveniente, dissidente	consciencioso, segue valores culturais, convencionais
Desenvoltura	acanhado, tímido, sensível	desenvolto, venturoso, insensível a repressões
Brandura	prático, objetivo, realista	sensível, harmonioso, sentimental
Confiança	confiante, acredita nas pessoas	desconfiado, suspeito, cauteloso

Fatores	Valores Baixos	Valores Altos
Imaginação	prático, cuidadoso, preciso, formal	imaginoso, regulado pelas solicitações interiores
Requinte	genuíno, sincero, simples	requintado, esmerado, isolado
Apreensão	plácido, seguro de si, sereno, complacente	apreensivo, indeciso, perturbado
Abertura a novas experiências	conservador, tradicional, dedicado à família	experimentador, renovador, liberal
Autossuficiência	dependente do grupo	autossuficiente, solitário, individualista
Disciplina	sem disciplina, tolerante a desordem, flexível	controlado, perfeccionista, organizado, autodisciplinado
Tensão	fleumático, relaxado, paciente	Tenso, impulsivo, impaciente

A tabela demonstra que para valores altos do fator “confiança”, por exemplo, tende-se a ser confiante, já para valores baixos o mesmo fator caracteriza como desconfiado e cauteloso. Anos mais tarde Cattell desenvolveu um questionário chamado "Questionário de 16 fatores de Personalidade". O questionário sofreu diversas revisões com o intuito de adaptar o instrumento a diferentes contextos, aplicações e culturas. Atualmente é composto por 185 questões de múltipla escolha, com tempo de aplicação de aproximadamente 50 minutos (NETO, 2010).

2.1.3. Modelo dos Três “Superfatores” de Eysenck

Hans Jurgen Eysenck foi um psicólogo influente no campo da Psicologia científica, e um dos pioneiros no estudo da estrutura fatorial da personalidade. A técnica fatorial determina quais comportamentos estão relacionados e são independentes de outros. Eysenck desenvolveu o modelo chamado "Modelo dos Três Superfatores" em bases biológicas dos traços, diferente do modelo apresentado anteriormente, cujo ponto de vista é baseado no léxico da linguagem. Em sua pesquisa Eysenck considera que uma dimensão da personalidade não é um traço de temperamento básico se não dispõe de uma base biológica apurada por meio de estudos correlacionais e experimentais. O modelo de Eysenck é considerado uma verdadeira teoria da personalidade, pelo motivo de apresentar uma validação experimental das propriedades dos traços (NETO, 2010).

O modelo de Eysenck inclui três dimensões tipológicas básicas, hierárquicas sendo, no nível inferior encontradas as respostas específicas e, no nível superior, encontrados os tipos, que são grupos de características estáveis e recorrentes do indivíduo. Diante disso, o modelo é

composto por três tipos estruturais da personalidade humana, que são (SANTOS; FLORES-MENDOZA, 2010):

- Extroversão: valores altos neste fator descrevem o sujeito como sociável, animado, ativo, assertivo, despreocupado, dominante, cordial, aventureiro e com busca de sensações. Valores baixos reúnem características opostas e definem o sujeito como introvertido.
- Neuroticismo: valores altos neste fator são definidos como ansiosos, deprimidos, tensos, irracionais, tímidos, melancólicos, emotivos, com tendência a sentir culpa e baixa autoestima. Valores baixos neste fator caracterizam o sujeito como emocionalmente estável.
- Psicoticismo: descrito, quando há altas pontuações, por adjetivos como agressivo, frio, egocêntrico, impessoal, impulsivo, antissocial, não-empático, criativo e obstinado. Valores baixos neste fator apresentam características contrárias e são definidos pelo controle de impulsos.

Para medir os valores de cada um dos fatores, Eysenck desenvolveu diversos questionários, entre eles: *Maudsley Personality Inventory*, *Eysenck Personality Inventory* e *Eysenck Personality Questionnaire*. Esses questionários, semelhantes ao de Cattell, foram revisados e adaptados inúmeras vezes.

As dimensões Neuroticismo e Extroversão são as mesmas apresentadas no modelo Cinco Grandes Fatores de personalidade, descrito na próxima seção.

2.1.4. Modelo dos Cinco Grandes Fatores

O modelo dos Cinco Grandes Fatores, ou *Big Five*, pode ser considerado uma versão contemporânea da teoria de traços que representa um ganho conceitual e empírico no estudo da personalidade, descrevendo dimensões humanas de forma consciente e replicável (HUTZ et al., 1998). Planejado na década de 1930, pelo psicólogo William McDougall (MCDUGALL, 1932), o modelo começou a ganhar expressão a partir da década de 1980, quando as pesquisas começaram a comprovar a existência de cinco traços básicos de personalidade em indivíduos de diferentes culturas e faixas etárias.

Utilizando uma abordagem léxica em linguagem natural, foi desenvolvido o modelo que descreve e classifica a personalidade humana em cinco fatores, sendo eles: Extroversão, Neuroticismo, Socialização, Realização e Abertura à experiência. Todos os fatores possuem

uma variedade de traços psicológicos, e o método utilizado para a escolha de cada fator parte do pressuposto que todos os descritores da personalidade humana que tem alguma relevância, interesse e importância estão registrados na linguagem natural (GARCIA, 2006).

Os cinco fatores que o modelo adota como constitutivo da personalidade das pessoas são descritos a seguir (MCCRAE; JOHN, 1992):

- Extroversão: caracteriza uma pessoa sensível, assertiva, ativa, impulsiva, sociável e que expressa entusiasmo;
- Neuroticismo: descreve uma pessoa insegura, ansiosa, mal-humoradas, autopunitivas e dimensões do afeto negativo;
- Socialização: descreve um sujeito como amigável, cooperativo, cordial, prestativo, altruísta e confiante;
- Realização: caracteriza uma pessoa auto-disciplinada, organizada, metódica e persistente;
- Abertura à experiência: descreve uma pessoa com abertura ao novo, intelectual, criterioso, liberal e tolerante.

No Brasil, os cinco fatores têm sido chamados conforme descrição supracitada, embora a literatura internacional tenha apontado algumas divergências em relação aos nomes (URQUIJO, 2001). Apesar de existir discrepâncias na forma como são chamados alguns fatores, as definições são consensuais e apontam para características semelhantes. Dessa maneira, será utilizada a nomenclatura mencionada nos estudos de (SILVA; NAKANO, 2011; NUNES, 2012). Para melhor entendimento de cada um dos fatores, a Tabela 2.2 apresenta alguns adjetivos que os caracterizam.

Tabela 2.2: Adjetivos característicos dos Cinco Grandes Fatores (SOUZA; NUNES, 2011).

	Extroversão	Socialização	Realização	Neuroticismo	Abertura
Polo do Rótulo	Ativo	Altruísta	Confiável	Ansioso	Artístico
	Aventureiro	Amigável	Consciente	Apreensivo	Curioso
	Barulhento	Carinhoso	Eficiente	Emotivo	Engenhoso
	Energético	Confiante	Minucioso	Instável	Esperto
	Entusiástico	Cooperativo	Organizado	Nervoso	Imaginativo
	Exibido	Gentil	Prático	Preocupado	Inteligente
	Sociável	Sensível	Preciso	Temeroso	Original
	Tagarela	Simpático	Responsável	Tenso	Sofisticado
Polo Oposto	Acanhado	Antipático	Desatento	Calmo	Comum
	Introvertido	Brigão	Descuidado	Contido	Simple
	Quieto	Bruto	Desorganizado	Estável	Superficial
	Reservado	Crítico	Distraído	Indiferente	Tolo
	Silencioso	Frio	Imprudente	Sereno	Trivial
	Tímido	Insensível	Irresponsável	Tranquilo	Vulgar

Segundo o modelo dos cinco Grandes Fatores, para uma classificação da personalidade é necessário conhecer a pessoa nos cinco traços. Para isso, alguns questionários foram criados com o intuito dessa avaliação. Destaca-se como principal o questionário NEO-PI-R (*Revised NEO Personality Inventory*), criado em 1992 pelos pesquisadores Paul Costa e Robert McCrae, que dividiram cada um dos cinco fatores em seis facetas inter-relacionadas, apresentadas na Tabela 2.3 (NETO, 2010).

Tabela 2.3: Facetas do questionário NEO-PI-R, adaptado de (NETO, 2010).

Extroversão	Socialização	Realização	Neuroticismo	Abertura
Acolhimento	Confiança	Competência	Ansiedade	Fantasia
Gregarismo	Franqueza	Ordem	Hospitalidade	Estética
Assertividade	Altruísmo	Senso de dever	Depressão	Sentimentos
Atividade	Aquiescência	Direcionamento	Autoconsciência	Ações
Busca de sensações	Modéstia	Autodisciplina	Impulsividade	Ideias
Emoções positivas	Sensibilidade	Deliberação	Vulnerabilidade	Valores

A seguir serão apresentados os principais testes psicológicos baseados no modelo *BigFive*.

2.1.5. Inventários de personalidade do Modelo *BigFive*

A maneira tradicional que psicólogos extraem os traços humanos – baseados nos cinco traços de personalidade do modelo *BigFive* – é geralmente por meio de instrumentos de mensuração de personalidade, como por exemplo, os testes psicológicos também chamados de inventários.

Segundo Pasquali (2001) um teste psicológico pode ser definido como um conjunto de tarefas determinadas antecipadamente, que um indivíduo precisa executar em uma situação normalmente artificializada ou sistematizada, tendo seu comportamento observado, descrito e julgado, sendo que essa descrição geralmente é feita utilizando-se números.

Os inventários apoiados em estudos empíricos são capazes de descrever as diferenças psicológicas entre as pessoas (NUNES, 2012). Como referido anteriormente, a abordagem dos traços que consiste em identificar características que melhor descrevem uma pessoa e quanto elas diferem das outras, influenciou o desenvolvimento de inventários para inferir a personalidade. Dessa maneira, existem vários testes psicológicos baseados no modelo *BigFive*: TIPI (*Ten-Item Personality Inventory*) (GOSLING et al., 2003); FFPI (*Five Factor Personality Inventory*) (HENDRIKS, 1997);); NEO-IPIP 120 (*Neo-International Personality*

Item Pool) (JOHNSON, 2014); BFQ (*Big Five Questionnaire*) (CAPRARA, et al., 1993); NEO-PI-R (*Revised NEO Personality Inventory*) (MCCRAE; COSTA, 1999); NEO-IPIP 300 (*Neo-International Personality Item Pool*) (JOHNSON, 2000); GPI (*Global Personality Inventory*) (SCHIMIT et al., 2002); dentre outros.

Para cada um dos inventários citados acima, as quantidades de questões são diferentes entre si. A quantidade de perguntas acaba influenciando a precisão das características mensuradas, isto é, quanto maior a quantidade de itens a serem avaliados, mais precisas serão as extrações do traço. Todavia, cada um dos instrumentos possuem particularidades e influências na metodologia de extração da personalidade. Para elucidar o tema, a Tabela 2.4 apresenta o número de questões relacionadas a cada inventário citado.

Tabela 2.4: Quantidade de questões dos inventários baseados no modelo *BigFive*.

Inventário	Número de Questões
TIPI (<i>Ten-Item Personality Inventory</i>)	10
FFPI (<i>Five Factor Personality Inventory</i>)	100
NEO-IPIP 120 (<i>Neo-International Personality Item Pool</i>)	120
BFQ (<i>Big Five Questionnaire</i>)	132
NEO-PI-R (<i>Revised NEO Personality Inventory</i>)	240
NEO-IPIP 300(<i>Neo-International Personality Item Pool</i>)	300
GPI (<i>Global Personality Inventory</i>)	504

Destaca-se com maior número de questões o inventário GPI (*Global Personality Inventory*). De acordo com DeRaad e Perugini (2002) o GPI é o maior inventário de personalidade baseado no modelo *BigFive*, contemplando 504 itens e categorizando 32 facetas de personalidade. Entretanto, a aplicação desse tipo de questionário em um ambiente computacional muitas vezes se torna inviável, devido ao grande número de questões a serem preenchidas pelo usuário e o tempo necessário para o preenchimento, o que pode levar a desistência da conclusão do teste, deixando o processo de avaliação do comportamento vulnerável.

Em decorrência disso, uma alternativa encontrada pelos pesquisadores (NUNES; TELES; DE SOUZA, 2013; QIU et al., 2012) que utilizam o preenchimento de inventários de maneira computacional e *online*, é a utilização de inventários com menos questões sem perder

a precisão nas representações dos traços de personalidade, como por exemplo NEO-IPIP 300 (*Neo-International Personality Item Pool*).

O inventário NEO-IPIP 300 foi criado e validado por Johnson (2000), com o objetivo de gerar uma versão gratuita do inventário *Neo Personality Inventory* (NEO-PI-R), o qual é descrito como um dos inventários comerciais mais robustos, conhecidos e validados no âmbito científico (JOHNSON, 2000). As questões do NEO-IPIP 300 estão divididas uniformemente entre cinco fatores do *BigFive*, cada fator é representado por um conjunto de 60 perguntas, totalizando 300 questões. Ainda, para cada fator são fornecidas seis facetas de personalidade, que representam dez questões para cada faceta. Todas as perguntas do NEO-IPIP 300 estão pontuadas em uma escala numérica de 1 a 5, que é associada à resposta do usuário. A Tabela 2.5, ilustra uma amostra de questões do inventário NEO-IPIP 300 e a relação com os fatores.

Tabela 2.5: Questões do inventário NEO-IPIP 300 baseados no modelo *BigFive*, adaptado de (NUNES, 2008).

Número da Questão	Fator do <i>BigFive</i> ¹	Faceta	Questão
1	N1	Ansiedade	Preocupo-me com as coisas
2	E1	Amigabilidade	Faço amigos facilmente
3	A1	Imaginação	Tenho uma imaginação vívida
4	S1	Confiança	Confio nos outros
5	R1	Auto-eficácia	Completo tarefas com sucesso
6	N2	Raiva	Fico com raiva facilmente
7	E2	Gregarismo	Adoro festas com muitas pessoas
8	A2	Interesses artísticos	Acredito na importância da arte
9	S2	Moralidade	Nunca sonegaria impostos
10	R2	Ordem	Gosto de ordem

Recentemente, Johnson (2014) desenvolveu uma nova versão do NEO-IPIP 300 com exatamente as mesmas características do original (30 facetas), mas de forma mais eficiente, com menos itens no questionário, cerca de 120 questões, sendo intitulado de NEO-IPIP 120. Segundo o autor, no inventário NEO-IPIP contendo 300 questões, a maioria das pessoas

¹ Dimensões do *BigFive* (N = Neuroticismo; E = Extroversão; S = Socialização; R = Realização; A = Abertura à experiência)

demora entre 20 a 40 minutos para a conclusão do questionário, para a nova versão de 120 questões o tempo de preenchimento é de 10 a 20 minutos.

A curta versão do NEO-IPIP fornece uma alternativa para as pessoas que não possuem tempo para concluir o inventário original. Embora a versão mais longa possua maior confiabilidade, por questões óbvias em relação ao tamanho do questionário, a versão curta, segundo Johnson (2014), atende aos padrões profissionais de confiabilidade e possui segurança de medição aceitável.

O modelo *Big Five* nos últimos anos tem sido um amplo suporte no âmbito científico, demonstrando que método tem replicabilidade a partir de diferentes sistemas teóricos, culturas e línguas conforme apresentado na pesquisa de (BARRICK et al., 2001). Para McCrae e Costa (MCCRAE; COSTA, 1999) a generalização da estrutura dos Cinco Fatores é uma classificação muito adequada e frequentemente confiável, ainda que não haja uma concordância universal.

Diante disto, a comunidade científica da área da Computação Afetiva (ramo da Inteligência Artificial) tem utilizado métodos e teorias da personalidade para que máquinas e sistemas computacionais reconheçam e classifiquem a personalidade de pessoas, visto que o uso do computador para comunicação e interação entre os indivíduos têm sido intensos nos últimos anos. Na próxima seção serão aduzidos os grandes diferenciais da área para tal inferência.

2.2. Inferência de Personalidade por meio de Textos

A identificação e inferência de personalidade humana têm sido relacionadas à área da Computação Afetiva (NUNES, 2012; NETO, 2010), com raízes na Inteligência Artificial. Nesta área estuda-se como os computadores podem discernir, modelar e responder às emoções humanas e, dessa maneira, como podem expressá-las por meio de uma interface ou interação computacional (PICARD, 1997).

Ainda, além da inferência de personalidade, a Computação Afetiva está relacionada ao reconhecimento de emoção, afeto, sentimento e opinião. Tais termos, segundo (MUNEZERO et al., 2014) possuem divergências e reúnem diversos conceitos, como: (i) afeto é o mais abstrato de ser verificado em linguagem, pois ele é considerado não consciente para o ser humano e é predecessor para sentimentos e emoções; (ii) sentimentos são fenômenos conscientes, centrados na pessoa; (iii) emoções são expressões sociais dos sentimentos

influenciadas pela cultura; e (iv) opiniões são interpretações pessoais de informações de determinada entidade que podem ou não conter sentimentos ou emoções. Para mais informações sobre as diferenças dos termos, sugere-se a leitura de (MUNEZERO et al., 2014).

Dessa maneira, a Computação Afetiva possui grande difusão nos estudos e pesquisas envolvendo a teoria da personalidade, devido ao seu potencial de aplicabilidade, sendo utilizada em diferentes áreas: (i) na Inteligência Artificial para criar agentes mais humanos (DIMURO, 2007), (ii) mensurar a personalidade dos usuários em redes sociais (QUERCIA et al., 2012; MAKOVIKJ et al., 2013; BACHRACH et al., 2012; TOMLINSON; HINOTE; BRACEWELL, 2013; NUNES, 2013), (iii) personalização de produtos em sites de negócios (e-commerce) (NUNES; CERRI; BLANC, 2008), e (iv) sistemas de recomendação baseados em personalidade (HU; PU, 2009).

Todos os trabalhos citados anteriormente têm como ênfase a inferência da personalidade do usuário por meio de texto escrito por ele. Essa abordagem será explorada na presente pesquisa, entretanto, os estudos de personalidade na área da computação afetiva contemplam outras fontes, tais como a linguagem falada (MAIRESSE, 2007), fotos em mídias sociais (CELLI et al., 2014; GUNTUKU et al., 2015), vídeos (BIEL; GATICA-PEREZ, 2012) e músicas (FERWERDA et al., 2015), as quais não serão abordadas.

Os trabalhos sobre o reconhecimento automático de personalidade são relativamente recentes. O estudo pioneiro abordando essa questão foi de Argamon, (ARGAMON et al., 2005) em 2005. O estudo se concentrou na descoberta de apenas dois traços de personalidade do modelo *BigFive*, sendo eles Extroversão e Neuroticismo. Adotou-se o método de identificação do estilo que os autores empregam em seus textos, partindo da noção intuitiva de que o estilo é indicado por características que representam a escolha do autor em utilizar uma expressão para um determinado conteúdo. Para isso, os autores aplicaram uma metodologia para a construção de um léxico usando como base os princípios da Gramática Sistêmico-Funcional².

Da mesma forma Oberlander e Nowson, em 2006, (OBERLANDER; NOWSON, 2006) trabalharam na classificação automática de personalidade a fim de detectar quatro dos cinco grandes traços (Extroversão, Neuroticismo, Socialização e Realização) em um *corpus* de *blogs* pessoais, considerados como uma espécie de “diário da internet”. Eles exploraram os

² Gramática Sistêmico-Funcional é uma abordagem funcional para análise linguística, criada por Halliday (HALLIDAY, 1994).

classificadores Naive Bayes e *Support Vector Machine* (SVM), treinando os conjuntos com diferentes recursos de n-gram.

Por sua vez, Mairesse (MAIRESSE et al., 2007) trabalhou no reconhecimento de todos os cinco grande traços de personalidade, tanto para texto quanto para conversa. O autor aludiu uma longa lista de correlações entre traços de personalidade do *BigFive* e dois conjuntos de léxicos: *Linguistic Inquiry and Word Count* (LIWC) e *Medical Research Council* (MRC), sendo que o primeiro inclui palavras classificadas como "emoções positivas" ou "raiva" e o último inclui estatísticas de palavras, tais como: estimativas de idade, frequência de uso e familiaridade. Maiores detalhes sobre LIWC e MRC, serão explorados na próxima seção.

A abordagem de Mairesse pode ser resumida em cinco passos: (i) coleta de *corpus* de forma individual; (ii) coleta de avaliações de personalidade para cada participante; (iii) extração de características importantes dos textos; (iv) construção de modelos estáticos de personalidade baseado em recursos de classificações; e (v) testes dos modelos de aprendizagem nas saídas linguísticas dos indivíduos.

Foram utilizados dois tipos de *corpus* para o estudo. O primeiro contém 2.479 redações (1,9 milhões de palavras) de estudantes de Psicologia norte-americanos, que foram orientados a escrever o que vêm a sua mente durante 20 minutos. Os dados foram recolhidos e analisados por Pennebaker e King (1999) e chamado de “*Essays*”. A personalidade desse *corpus* foi avaliada através do preenchimento de cada participante ao questionário *Inventory Big Five* de (JOHN et al., 1991).

A segunda fonte de dados é constituída por conversas gravadas usando um gravador eletrônico, recolhido por (MEHL et al., 2006) denominado de “EAR”. Para preservar a privacidade dos participantes dessa gravação, foram registrados apenas trechos aleatórios de conversas. Este *corpus* é muito menor do que ao anterior, contendo 96 participantes para um total de 97.468 palavras e 15.269 declarações. Enquanto o *corpus Essays* consiste apenas de textos, o EAR contém ambos os extratos de som e transcrições. Além dos léxicos citados anteriormente, foram acrescentados para o *corpus* EAR métricas pertinentes a fala, como tipos de enunciados (por exemplo: verbos de comando) e prosódica.

Para a execução dos experimentos com a intenção de avaliar o modelo criado, o autor utilizou algoritmos de aprendizagem de máquina de classificação, regressão e *ranking*. Os algoritmos de classificação analisados foram: Árvore de decisão (J48), *Nearest neighbour* (NB), *Naive Bayes* (NB), Ripper (JRip), AdaBoost e Máquinas de Vetores de Suporte com

kernel linear (SMO). Para a regressão, foram utilizados cinco algoritmos que retornam a pontuação média de personalidade: *Linear Regression* (LR); Árvore de regressão M5 (M5R), Árvore M5 com um modelo linear (M5), REPTree (REP) e Máquinas de Vetores de Suporte para regressão (SMO). Em relação ao problema de *ranking*, foi usado *RankBoost*, um algoritmo para organizar um *ranking* a partir das entradas (FREUND et al., 1998).

A Tabela 2.6 ilustra os melhores resultados obtidos por Mairesse para cada traço de personalidade. O autor confrontou diversas bases de treinamentos, adicionando e removendo recursos e léxicos a fim de produzir o melhor modelo de inferência de personalidade. Cada linha da tabela contém o algoritmo, o conjunto de recurso e o desempenho do modelo.

Tabela 2.6: Comparação dos melhores modelos para cada traço. (1) Precisão da classificação; (2) porcentagem de melhoria em relação à *baseline* utilizando regressão; (3) perda de ranking.

Adaptado de (MAIRESSE et al., 2007).

Traços de Personalidade	Classificação			Regressão			Ranking		
			(1)			(2)			(3)
Baseline	n/a	none	50%	n/a	none	0%	n/a	none	0.50
Modelos treinados em dados escritos (essays)									
Extroversão	ADA	LIWC	56%	LR	MRC	1%	Rank	LIWC	0.44
Neuroticismo	SMO	LIWC	58%	M5	LIWC	4%	Rank	LIWC	0.42
Socialização	SMO	LIWC	56%	LR	LIWC	2%	Rank	LIWC	0.46
Realização	SMO	LIWC	56%	M5	LIWC	2%	Rank	LIWC	0.44
Abertura à experiência	SMO	LIWC	63%	M5	Todos	7%	Rank	LIWC	0.39
Modelos treinados em dados de voz (EAR)									
Extroversão	NB	Todos	73%	REP	LIWC	24%	Rank	Prosódia	0.26
Neuroticismo	NB	Todos	74%	M5	Prosódia	15%	Rank	MRC	0.39
Socialização	NB	Todos	61%	M5R	Todos	3%	Rank	Todos	0.31
Realização	NB	Todos	68%	M5R	LIWC	18%	Rank	Todos	0.33
Abertura à experiência	NB	Todos	65%	M5	Tp. Enun.	1%	Rank	LIWC	0.37

O estudo de inferência de traços de personalidade humana com a computação tem propiciado o avanço da personalização de informações, produtos e serviços aos clientes na rede mundial de computadores, a *Internet*, e tem influenciado usuários na tomada de decisão como, por exemplo, em ambientes educacionais (PORTO et al., 2011). Com o advento da *Web 2.0* e suas plataformas de *blogs*, fóruns de discussão e vários tipos de mídias sociais, grandes volumes de dados foram gerados contendo informações de como os usuários da rede

interagem entre si, expressões de pensamentos e opiniões, tornando-se possível computacionalmente processar esses dados, com o propósito de reconhecer a personalidade e emoções dos autores, permitindo ainda antecipar possíveis reações de comportamento. Com isso, o RPT tem se tornado o grande diferencial para tal inferência, sendo adotado em vários trabalhos e recebido crescente atenção nos últimos anos.

Nesse contexto, Celli (CELLI, 2011) propõe um modelo de reconhecimento de personalidade por meio de publicações do popular site de *micro-blogging* *Twitter*, coletando um total de 25.700 *post*. O autor utilizou uma lista de recursos linguísticos fornecida por Mairesse (MAIRESSE, 2007) que se correlaciona com os traços de personalidade do modelo *BigFive* para a língua inglesa. Ao todo são 12 recursos que auxiliaram a ferramenta de reconhecimento da personalidade, a saber (CELLI, 2011):

- Pontuação: Total de “., ; ” : encontrados na sentença;
- Vírgulas: Total de “,” encontradas na sentença;
- Referência a outros usuários do *Twitter*: Total de “@...” encontrados na sentença;
- Ponto de exclamação: Total de “!” encontrados na sentença;
- *Emoticons* negativos: Total de *emoticons* que expressam sentimentos negativos na sentença;
- Números: Total de números encontrados na sentença;
- Parênteses: Número total de frases entre parênteses na sentença;
- *Emoticons* positivos: Total de *emoticons* que expressam sentimentos positivos na sentença;
- Ponto de interrogação: Total de “?” encontrados na sentença;
- Palavras longas: número total de palavras com mais de 6 caracteres na sentença;
- Type/token ratio (tt): Definido na Equação 1.1, onde: (w) contagem de palavras já utilizadas na sentença; (T) contagem de palavras total na sentença;
- Número de palavras: Total de palavras na sentença.

$$tt = \frac{w - T}{T} \quad (1.1)$$

A ferramenta de Celli sobre reconhecimento de personalidade não precisa de dados anotados a fim de estabelecer o modelo da personalidade dos usuários. O sistema pode avaliar

a personalidade apenas para os usuários que têm mais de um *post*, enquanto os demais usuários (e seus respectivos *posts*) são descartados. Segundo o autor, o estudo foi pioneiro em executar uma ferramenta de reconhecimento de personalidade no *Twitter*.

Na mesma rede social, *Twitter*, Celli realizou outras pesquisas na área de reconhecimento de personalidade. Em (CELLI; ZAGA, 2013), o autor aludiu uma breve relação entre sentimentos e personalidade em publicações na rede social. O estudo utilizou a mesma ferramenta de reconhecimento citada acima, entretanto, acrescentou ao processo a seleção de atributos.

A seleção de atributos tem como objetivo a melhoria de um método de aprendizagem de máquina quanto a sua taxa de precisão, selecionando um subconjunto menor de atributos capaz de descrever o conceito alvo o mais próximo possível da utilização de todos os atributos. Para a tarefa de seleção, dois objetos devem ser configurados: um seletor de atributo e um método de pesquisa.

Nesse contexto, Celli utilizou o seletor de atributos chamado “CfsSubsetEval” (Correlation-based Feature Selection) (HALL, 1998), e como método de pesquisa o algoritmo “Greedy Stepwise”. O papel do seletor determina qual método é utilizado para atribuir um valor a cada subconjunto de atributos, enquanto o método de busca determina o tipo de pesquisa que será realizada dentro do espaço de combinações de atributos.

O seletor *CfsSubsetEval* utiliza uma correlação em subconjuntos, avaliando a capacidade de predição de cada atributo no subconjunto juntamente com o grau de redundância entre os atributos. É considerado bom o subconjunto cujos atributos são altamente correlacionados com a classe e que contém atributos não correlacionados entre si (PICCHI NETTO, 2013). Por sua vez, o *Greedy Stepwise* é utilizado para efetuar a pesquisa de características por meio de uma busca gulosa nos subconjuntos de cada espaço de atributos (WITTEN; FRANK, 2005). Sobre *CfsSubsetEval* e *Greedy Stepwise* sugere-se a leitura de (HALL, 1998) e (WITTEN; FRANK, 2005) respectivamente.

A seleção do conjunto de atributos, mencionada anteriormente, é um problema que vem sendo estudado há décadas, conforme comentado em (AHA; BANKERT, 1996). A partir dessas pesquisas foram produzidos inúmeros seletores de atributos. Os pesquisadores de RPT (POHJALAINEN; RÄSÄNEN; KADIOGLU, 2015) e (BENCHERIF et al., 2012) utilizaram-se do seletor “*Wrapper*”, proposto por Kohavi e John (1997). Essa abordagem avalia cada subconjunto de atributos por meio de um algoritmo de aprendizado, utilizando-se da acurácia

desse algoritmo como medida para a avaliação (PICCHI NETTO, 2013). Para mais detalhes sobre seleção de atributos, os estudos de (HALL; SMITH, 1998) e (KHALID; KHALIL; NASREEN, 2014) são sugeridos ao leitor.

Cumprе mencionar que além da pesquisa de Celli, vários outros estudos (QUERCIA et al., 2012; BACHRACH et al., 2012; SAIDMAN, 2012; MAKOVIKJ et al., 2013; KOSINSKI et al, 2013) emergem dos dados gerados pelas redes sociais, com intuito de classificar a personalidade e o impacto de seus autores.

Os estudos expostos anteriormente, utilizaram o modelo *BigFive*, que frequentemente é usado para esse tipo de análise, já que, quando comparado com outro modelo, ele descreve de forma clara e objetiva os traços de personalidade. Ainda, leva-se em consideração que o modelo possui grande aceitabilidade não apenas na área da computação afetiva, mas também se relaciona com um vasto número de estudos no âmbito científico.

Ressalta-se que todos os estudos citados possuem métodos de reconhecimento de personalidade a partir de textos escritos em língua inglesa. Grande parte da comunidade científica tem buscado a expansão e o melhoramento de métodos e léxicos para o auxílio de detecção de personalidade para o inglês. Consequentemente, foram construídos vários léxicos, a fim de incrementar a estrutura de métodos eficazes para a inferência de personalidade e emoções empregadas em textos escritos no idioma inglês. Os principais léxicos serão apresentados nas próximas seções.

2.2.1. Léxicos

Conforme Specia e Nunes (2004), os léxicos computacionais são recursos lexicais criados, geralmente, de forma manual, especificamente para o tratamento computacional. Em seu âmbito geral fornecem diversas características morfológicas e orientação semântica de uma palavra, tais como a identificação de classes gramaticas, negação, conjugação verbal e outros. Esta seção fornece uma descrição dos principais léxicos psicolinguísticos investigados para a pesquisa. Tais léxicos são os mais populares na literatura, ou seja, os mais citados e amplamente utilizados.

Linguistic Inquiry and Word Count (LIWC). Criado por (PENNEBAKER et al., 2001), é uma ferramenta de análise de texto que avalia componentes emocionais, cognitivos e estruturais de um determinado texto, baseia-se na utilização de um dicionário contendo classificação de palavras em categorias. Em sua atual versão (PENNEBAKER et al., 2007), o

LIWC fornece mais de 70 categorias, divididas em 4 dimensões, a saber: (i) linguísticas, contém as categorias de artigos, preposições, pronomes de primeira pessoa do singular, pronomes de primeira pessoa do plural, advérbios, negações, palavrões, entre outras; (ii) psicológicos, possui as categorias de emoções positivas e negativas, percepção cognitiva, entre outras; (iii) relatividade, com as categorias de tempo, tempo verbal, espaço e etc; e (iv) pessoal, contém as categorias sobre casa, trabalho, lazer, entre outras. A título de exemplo, a palavra “agree” no dicionário pertence a 5 categorias: *assent*, *affective*, *positive emotion*, *positive feeling*, e *cognitive process*. Portanto, o LIWC pode detectar mais de uma categoria para a mesma palavra. As categorias e dimensões do LIWC estão resumidas na Tabela 2.7.

Tabela 2.7: Características textuais do LIWC, adaptado de (PENNEBAKER et al., 2007).

Dimensões Categorias	Exemplos de palavras (Inglês)
Linguística	
Contagem de palavras	
Palavras por frase	
Palavras > 6 letras	
Total de pronomes	I, them, itself
Pronome pessoal	I, them, her
1ª pessoa do singular	I, me, mine
1ª pessoa do plural	We, us, our
2ª pessoa	You, your, thou
3ª pessoa do singular	She, her, him
3ª pessoa do plural	They, their, they'd
Pronomes impessoais	It, it's, those
Artigos	A, an, the
Verbos comuns	Walk, went, see
Verbos auxiliares	Am, will, have
Passado	Went, ran, had
Presente	Is, does, hear
Futuro	Will, gonna
Advérbios	Very, really, quickly
Preposição	To, with, above
Conjunções	And, but, whereas
Negação	No, not, never
Quantificadores	Few, many, much
Números	Second, thousand
Palavrões	Damn, piss, fuck
Psicológica	
Social	Mate, talk, they, child

Dimensões Categorias	Exemplos de palavras (Inglês)
Família	Daughter, husband, aunt
Amigos	Buddy, friend, neighbor
Humano	Adult, baby, boy
Afetivo	Happy, cried, abandon
Emoções positivas	Love, nice, sweet
Emoções negativas	Hurt, ugly, nasty
Ansiedade	Worried, fearful, nervous
Raiva	Hate, kill, annoyed
Tristeza	Crying, grief, sad
Cognitivos	Cause, know, ought
Intuição	think, know, consider
Causa	because, effect, hence
Discordância	should, would, could
Tentativa	maybe, perhaps, guess
Certeza	always, never
Inibição	block, constrain, stop
Inclusivo	And, with, include
Exclusivo	But, without, exclude
Perceptivo	Observing, heard, feeling
Ver	View, saw, seen
Ouvir	Listen, hearing
Sentir	Feels, touch
Biológico	Eat, blood, pain
Corpo	Cheek, hands, spit
Saúde	Clinic, flu, pill
Sexual	Horny, love, incest
Ingestão	Dish, eat, pizza
Relatividade	Area, bend, exit, stop
Movimento	Arrive, car, go
Espaço	Down, in, thin
Tempo	End, until, season
Pessoal	
Trabalho	Job, majors, xerox
Conquista	Earn, hero, win
Lazer	Cook, chat, movie
Dinheiro	Audit, cash, owe
Religião	Altar, church, mosque
Morte	Bury, coffin, kill
Falada	
Concordância	Agree, OK, yes
Sem fluência	Er, hm, umm
Preenchimento	Blah, I mean, you know

O principal objetivo do LIWC é agrupar palavras em categorias que podem ser utilizados para analisar as características psicolinguísticas em textos.

Recentemente o LIWC foi disponibilizado para idioma português por meio do estudo de (BALAGE FILHO et al., 2013). Os autores efetuaram a tradução e adaptação de 127.162 palavras do português catalogadas em 64 categorias, dentre elas: palavras relacionada à família, palavras, verbos, palavras de negação e outros. Destaca-se a utilização do LIWC nos trabalhos de (MAIRESSE, 2007; CELLI, 2012; PORIA; GELBUKH; AGARWAL et al., 2013). A Tabela 2.8 apresenta todas as categorias do LIWC que possui palavras da língua portuguesa.

Tabela 2.8: Dimensões e categorias do LIWC para o português (BALAGE FILHO et al., 2013).

Dimensões	Categorias
Linguística	função de palavras, total de pronomes, pronome pessoal, 1ª pessoa do singular, 1ª pessoa do plural, 2ª pessoa, 3ª pessoa do singular, 3ª pessoa do plural, pronomes impessoais, artigos, verbos comuns, verbos auxiliares, passado, presente, futuro, advérbios, preposição, conjunções, negação, quantificadores, números, palavras
Psicológica	social, família, amigos, humano, afetivo, emoções positivas, emoções negativas, ansiedade, raiva, tristeza, cognitivos, intuição, causa, discordância, tentativa, certeza, inibição, inclusivo, exclusivo, perceptivo, ver, ouvir, sentir, biológico, corpo, saúde, sexual, ingestão, relatividade, movimento, espaço, tempo
Pessoal	trabalho, conquista, lazer, dinheiro, religião, morte
Falada	concordância, sem fluência, preenchimento

MRC Psycholinguistic DataBase. Criado por (COLTHEART, 1981), o MRC é uma base de dados que contém categorias psicolinguísticas diferentes para cada palavra. Em sua versão inicial a base de dados é composta por 98.538 palavras e fornecida como um serviço *online*, disponibilizando vários arquivos e programas para seu acesso. Para a sua segunda versão foram incrementadas 52.299 novas entradas, totalizando 150.837 palavras, assim construindo um dicionário com ênfase em computadores, livremente para fins de processamento de linguagem natural e tarefas de inteligência artificial. A estrutura geral das duas versões é semelhante, contendo 26 propriedades psicolinguísticas diferentes para as palavras (WILSON, 1988). Ainda, o MRC não inclui apenas a informação sintática, mas também dados psicológicos para as palavras.

Utilizou-se o MRC nos trabalhos de (GOLBECK et al., 2011; MAIRESSE, 2007; CELLI, 2012; PORIA; GELBUKH; AGARWAL et al., 2013).

WordNet. É uma base de dados composta por entradas, ou palavras, que sistematiza o conjunto dos verbos, substantivos, adjetivos e advérbios de um idioma. Teve origem na Universidade de Princeton (MILLER, 1995; FELLBAUM, 1998), no ano de 1984, como um experimento linguístico anotado de forma manual, com o intuito de ser uma base lexical para a língua inglesa, contendo em sua estrutura níveis morfológicos e semânticos. A sinonímia é a principal relação entre palavras da *WordNet*, sendo que os grupos de palavras sinônimas são organizadas em conjuntos denominados *synsets* (do inglês *synonyms sets*), em que cada *synset* representa um conceito cujo sentido é válido para todas as palavras do conjunto (SCARTON, 2013). Atualmente o léxico é composto por 152.059 palavras, entre verbos, substantivos, adjetivos e advérbios e 115.424 *synsets*. Uma versão do *WordNet* foi criada pelos lingüistas brasileiros para o português, denominada *WordNet.br*, sendo uma versão menor que o léxico original, contendo 18.500 *synsets* e 44.000 tipos de palavras (DIAS-DA-SILVA et al., 2008).

Observa-se a utilização do *WordNet* nos trabalhos de (BOUCHET et al., 2010; ARYA et al., 2012; DIAS-DA-SILVA et al., 2008; SANSONNET; BOUCHET, 2010).

VerbNet. É um léxico verbal para a língua inglesa, criado por Karin Kipper (KIPPER-SCHULER, 2005), que fornece informações sintáticas, semânticas e uma descrição completa dos verbos inspiradas no trabalho de Levin (1993), resultando em uma coleção estruturada de classes verbais e alterações sintáticas. Em sua primeira versão, a partir das classes de Levin, foram anotados 191 classes com cobertura para 4.656 verbos. Atualmente, com o acréscimo de outras classes de verbos, já existe anotação para cerca de 5.800 verbos, divididos em 274 classes (PALMER, 2005). Com o objetivo de gerar um léxico com as mesmas características para o português o Brasil, os pesquisadores Scarton e Aluísio (2012) criaram o *VerbNet.br*. Esse estudo considera a compatibilidade de tradução entre os verbos do *VerbNet* (*cross-linguístico*) para o português, herdando assim os recursos semânticos do *VerbNet*, essa transposição foi feita de maneira semiautomática.

2.2.2. Léxicos Afetivos

A partir do desenvolvimento de métodos, técnicas e recursos que, integrados, tornam os sistemas computacionais capazes de manipular significado afetivo e de sentimentos no

discurso, surgem os léxicos especializados em termos afetivos que dão suporte ao desenvolvimento de trabalhos nessas áreas (PASQUALOTTI; VIEIRA, 2008).

De acordo com Ortony e colegas (1987), um léxico afetivo permite utilizar não somente palavras que se referem diretamente a emoções, mas a muitas outras palavras, que mesmo não se referindo a emoções, implicam em diversas outras formas.

A abordagem de léxico afetivo está dividida em: (i) abordagem em dicionário e (ii) abordagem em *corpus*. A primeira abordagem possui um conjunto de termos emocionais recolhido de forma manual, em que cada termo possui a sua emoção. Na segunda abordagem o léxico é construído a partir de termos emocionais rotulados manualmente e novos termos emocionais são acrescentados com orientações específicas de contexto, utilizando-se métodos estatísticos ou semânticos (LIU, 2012). A Figura 2.1 mostra como as abordagens de léxicos estão subdivididas.

Na tarefa de reconhecimento de personalidade, em que léxicos afetivos são usados, a abordagem baseada em dicionário é comumente utilizada, devido ao fato de ser livre de orientações específicas de contexto. Dessa forma, o presente estudo visa explorar a utilização de léxicos baseados em dicionário, haja vista que atualmente para o português não há disponibilidade de léxicos baseados em *corpus*. Para detalhes da utilização da abordagem léxica baseada em *corpus*, sugere-se a leitura dos estudos (KAJI; KITSUREGAWA, 2006), (WIEBE; MIHALCEA, 2006), (KAJI; KITSUREGAWA, 2007) e (WU, 2010).

Em seguida, é dada uma visão geral dos principais léxicos afetivos (abordagem em dicionário) utilizados na literatura.

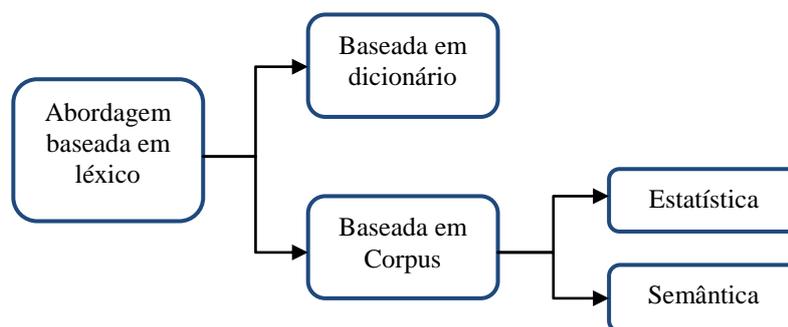


Figura 2.1: Divisão da abordagem baseado em léxico para a Análise de Sentimento em textos.

Adaptada de (MEDHAT; HASSAN; KORASHY, 2014).

SentiStrength. É um classificador baseado em léxico, construído por (THELWALL et al., 2010) com a finalidade de detectar sentimentos em textos curtos da língua inglesa. *SentiStrength* combina uma abordagem baseada em léxico com regras linguísticas mais sofisticadas, por exemplo, erros ortográficos, pontuação e uso de *emoticons*.

Para cada entrada de texto, o *SentiStrength* classifica as palavras em um intervalo de 1 a 5 denotando palavras positivas, e -5 a -1 para palavras negativas (GARCIA; SCHWEITZER, 2011). Por exemplo, um texto com uma pontuação de 3 e -5 contém sentimento positivo moderado e forte sentimento negativo, respectivamente. A frase “*I love you but hate the current political climate.*” é classificada da seguinte maneira: “*I love[3] you but hate[-4] the current political*”, para cada palavra avaliada da frase é extraída uma pontuação mínima de -5 e máxima de 5, pertinente ao seu conteúdo emocional inserido no léxico (SENTISTRENGTH, 2015).

Além de reportar a polaridade emocional (positivo e negativo) da palavra, pode-se ainda extrair a escala trinária (negativo, neutro e positivo) e escala única. Essa última escala reporta um único resultado para todo o texto, em um intervalo de 1 (não positivo) a 5 (extremamente positivo) e -1 (não negativo) a -5 (extremamente negativo) (THELWALL, 2013). Seguindo o exemplo da frase “*I love you but hate the current political climate.*”, o resultado de escala única é -1.

O *SentiStrength* possui um ferramenta *online* (<http://www.sentistrength.wlv.ac.uk/>) que permite a detecção de sentimento para textos curtos. Originalmente desenvolvido para inglês, permite a classificação de sentimento em outros idiomas, entre eles o português, porém, sem o nível de precisão existente na língua inglesa (SENTISTRENGTH, 2015).

A utilização do *SentiStrength* está presente nos estudos de (TENORIO, 2014; THELWALL, 2013; FARNADI et al., 2014).

SentiWordNet. É uma ferramenta resultante de um estudo realizado e mantido pelo *Istituto di Scienza e Tecnologie dell’Informazione “A. Faedo”*; pertencente ao Conselho Nacional de Pesquisa na Itália (ESULI; SEBASTIANI, 2006; SENTIWORDNET, 2015). Sua finalidade é a mineração de sentimentos utilizando o dicionário *WordNet* (MILLER, 1995) como fonte de dados. A relação entre os dicionários é feita por meio da associação de cada *synset* com três valores de pontuação, sendo eles positivo, negativo e neutro, que indicam o sentimento de um texto. Cada palavra identificada recebe uma pontuação variando de 0 a 1 (TENORIO, 2014). Para exemplificar o comportamento deste léxico, a Figura 2.2 ilustra a

submissão da palavra “good” para a classificação, por meio de um sistema *web*³. Observa-se que, dependendo do contexto, ela pode ser avaliada como uma emoção positiva ou neutra (Adjetivo) e usualmente é associada a sentimentos positivos para indicação de substantivo. Ressalta-se que o *SentiWordNet* foi desenvolvido exclusivamente para a língua inglesa, não havendo expansão para outros idiomas.

Exemplificando sua abrangência, podemos citar a utilização do *SentiWordNet* nos estudos de (DENECKE, 2008; CHIKERSAL et al., 2015; MOSTAFA, 2013).



Figura 2.2: Exemplo de classificação *SentiWordNet*. Adaptada de (SENTIWORDNET, 2015).

WordNetAffect. É uma extensão da base de dados *WordNet* (FELLBAUM, 1998), criado pelos pesquisadores (STRAPPARAVA; VALITUTTI, 2004), com o propósito de adquirir subconjunto de *synsets* da base *WordNet* adequados para representar conceitos afetivos. Para a construção desse léxico, primeiramente foi necessário criar uma base (*Affect*) composta pelas classes gramaticais da *Wordnet* e outras informações lexicais, semânticas e afetivas. As informações lexicais e semânticas referem-se: (i) correlação entre as línguas inglesa e italiana, (ii) classe gramatical à qual a palavra pertence, (iii) relações de sinonímia e antonímia, e (iv) um glossário com uma breve descrição (PASQUALOTTI; VIEIRA, 2008). As informações afetivas, por sua vez, referem-se às teorias de emoções baseadas no conceito da avaliação cognitiva (ORTONY et al., 1987), às teorias das emoções básicas (ELLIOT, 1992) e às teorias dimensionais, representando sua valência afetiva: positivo ou negativo. Após a base *Affect* criada, os autores realizam duas novas etapas para a construção do

³ <http://sentiwordnet.isti.cnr.it/index.php>

WordNetAffect (STRAPPARAVA; VALITUTTI, 2004). A primeira etapa consistiu na projeção dos termos da base *Affect* para os respectivos *synsets* na base *WordNetAffect*, com uma etiqueta afetiva, denominada de “*a-label*”, que possui o conteúdo da teoria cognitiva da base *Affect*.

A seguir, a Tabela 2.9 apresenta a etiqueta “*a-label*” com os respectivos estados afetivos e exemplos.

Tabela 2.9: Lista de “*a-labels*” com os respectivos estados afetivos e exemplos (PASQUALOTTI; VIEIRA, 2008).

a-label	Estado afetivo	Exemplos
Emo	Emotion	substantivo “anger”, verbo “fear”
Moo	Mood	substantivo “animosity”, adjetivo “amiable”
Tra	Trait	substantivo “aggressiveness”
Cog	cognitive state	substantivo “confusion”, adjetivo “dazed”
Phy	physical state	substantivo “illness”, adjetivo “all_in”
Eds	edonic signal	substantivo “hurt”, substantivo “suffering”
Sit	emotion- situation eliciting	substantivo “awkwardness”
Res	emotional response	substantivo “cold_sweat”, verbo “tremble”
Beh	behaviour	substantivo “offense”, adjetivo “inhibited”
Att	attitude	substantivo “intolerance”, noun “defensive”
Sen	sensation	substantivo “coldness”, verbo “feel”

Para a segunda etapa, foram realizadas as correlações com os *synsets* da *Wordnet*, em que as relações semânticas foram utilizadas para alimentar a *WordnetAffect*. A cada relação semântica da *WordNet* (antônimo, similaridade, “pertencer a”, atributos, “ver também” e “derivado-de”), foi verificado se o significado afetivo era preservado. As relações de hiperônimo, implicação, causas e grupos de verbos, em que o significado afetivo é parcialmente preservado, não foram considerados para compor a base afetiva proposta (PASQUALOTTI; VIEIRA, 2008). Para exemplificar, a Tabela 2.10 apresenta a etiqueta que contém informações afetivas (coluna “CAT”) para alguns registros da base *WordnetAffect*. Ainda, a tabela contém: a classe gramatical do *synset* (POS - *Part Of Speech*), código de localização do *synset* na base *WordNet* (ID do *Synset*) e a maneira como o *synset* foi obtido, podendo ser da base *Affect (core)* ou das relações da *WordNet (Origem)*.

Uma tradução e adaptação do *WordnetAffect* foi realizada para o português brasileiro, denominada *WordnetAffectBR* (PASQUALOTTI; VIEIRA, 2008). Os pesquisadores nortearam a pesquisa para que a tradução não contemplasse somente o sentido isolado das palavras, mas para considerar o seu significado presente nas definições da *Wordnet* e na estrutura dos *synsets*. Por isso, além das palavras, o glossário de cada *synset* também foi traduzido. Dessa forma, a base original da *WordnetAffectBR*, na fase final de tradução, apresentou 457 palavras após a retirada das duplicidades e análise dos significados, a base ficou composta por 289 palavras (PASQUALOTTI; VIEIRA, 2008).

Tabela 2.10: Estrutura da *WordnetAffect* e exemplos de registros da base (PASQUALOTTI; VIEIRA, 2008).

POS	ID do <i>Synset</i>	Origem	CAT
N	#10972097	core	emo
N	#03848510	core	beh, att, tra
V	#00548199	antonym	emo, cog
A	#02475653	pertains-to	eds,emo

O WordNetAffect está presente nos trabalhos de (PORIA; GELBUKH; HUSSAIN, 2013; VALITUTTI; STRAPPARAVA, 2010).

Affective Norms for English Words (ANEW). Pesquisadores do *National Institute of Mental Health* (NIMH), na Universidade da Flórida, desenvolveram diferentes conjuntos de estímulos emocionais e que fornecem classificações normativas de valência, alerta e dominância, são eles: *International Affective Picture System (IAPS)*, *International Affective Digital Sounds (IADS)*, *Affective Norms for English Text (ANET)* e *Affective Norms for English Words (ANEW)*. O ANEW (BRADLEY; LANG, 1999) traz o conceito que a emoção pode ser definida como uma relação temporal breve, composta pelo menos de duas dimensões, uma de valência (do desagradável ao agradável) e outra de alerta (do relaxado ao estimulado) (KRISTENSEN *et al.* 2011). Dessa forma é possível caracterizar um estímulo de uma palavra, medindo o quão desagradável ou agradável, referindo-se à valência e ao quão relaxado ou estimulado ficamos perante a palavra, referindo-se ao alerta.

Seu conjunto de palavras é composto por 1.034 vocábulos com medidas para três dimensões emocionais. A primeira dimensão diz respeito à valência, que consiste na avaliação do quão agradável ou desagradável é a percepção de uma palavra, isto é, a assimilação de um

indivíduo a um vocábulo que provoca uma emoção discreta de felicidade pode ser classificada por uma valência agradável. A segunda dimensão é chamada de alerta, que avalia o quão estimulada ou relaxada é a percepção de uma palavra, como por exemplo, palavras que evocam raiva podem ser classificadas com alerta alto. A última dimensão, chamada de dominância, consiste na análise do quão em controle de uma palavra ou dominado por ela nós a percebemos.

Considerando sua utilidade e relevância para o meio científico, o ANEW foi traduzido para uma versão em português, sofrendo adaptação e normalização para a população brasileira. A norma brasileira para o *affective norms for english works* (ANEW-Br) (KRISTENSEN *et al.* 2011) teve como objetivo obter medidas de valência e alerta para um conjunto de 1.046 palavras em português, realizando um estudo de tradução, adaptação e normatização do ANEW.

Para determinar o nível emocional do conjunto final de 1.046 palavras em português, foi utilizada a escala de *Self-Assessment Manikin* (SAM) (BRADLEY; LANG, 1994). Sendo essa escala formada por nove pontos, conforme ilustrado na Figura 2.3, em que cinco desses pontos são apontados por figuras de bonecos que indicam graus de reações emocionais, possuindo uma escala própria ilustrando variações para determinada dimensão. O SAM que representa valência na figura (1A) varia de um boneco sorridente a um boneco uma expressão descontente, para o SAM que representa o alerta, (1B) sua caracterização tem sua variação de um boneco ativo até um boneco inerte ou relaxado.

Para a elaboração do ANEW-Br, primeiramente foi realizada a tradução e adaptação das palavras do ANEW para o português brasileiro, traduzindo os 1.034 vocábulos de forma direta, removendo todas as traduções que resultaram em palavras compostas. Também, novas palavras emocionais foram adicionadas, gerando uma lista de 1.046 palavras para o idioma português brasileiro. Posteriormente em outra etapa foi realizado o julgamento de valência e alerta para as palavras resultantes (ANEW-Br).

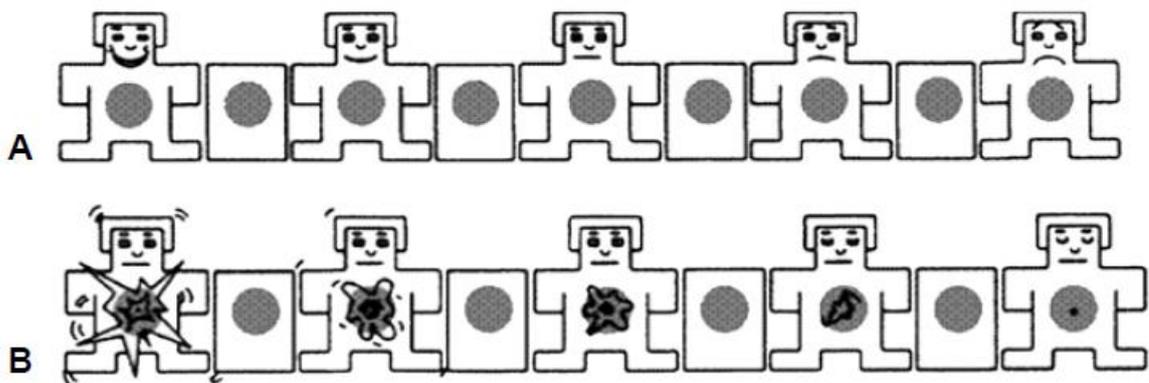


Figura 2.3: Escala de avaliação de valência (A) e alerta (B) do *Self-Assessment Manikin*.

(KRISTENSEN *et al.* 2011).

Por fim, é importante elucidar que as normas brasileiras e americanas possuem diferenças. Kristensen (KRISTENSEN *et al.* 2011) relata que o ANEW-Br não é apenas uma tradução e normatização do ANEW, possui adaptação das normas americanas, devido à inadequação de traduções do inglês para o português e vice-versa.

Exemplifica-se a abrangência do ANEW nos trabalhos de (NIELSEN, 2011; YOUNG; SOROKA, 2012).

SentiNet. É um léxico de proporciona a inferência de polaridade de textos em inglês em nível semântico, ao invés de sintático, isto é, para identificação de emoções não se baseia em partes do texto em que as opiniões são expressas de forma explícita, como termos positivos (por exemplo: bom, agradável, excelente, correto, superior, melhor) ou termos negativos (por exemplo: ruim, desagradável, infeliz, errado, inferior, pior) e sim por meio de informações de senso comum (*sentic computing*⁴), *Web Semântica*, e técnicas de computação afetiva (CAMBRIA *et al.*, 2010).

O *SenticNet* extrai sentimento e informação semântica de 30.000 conceitos de conhecimento do senso comum. Os conceitos são extraídos usando um método que retorna variáveis de sentimento associado a cada um dos conceitos encontrados em uma mensagem, como: a pontuação de polaridade e o vetor *Sentic* (CAMBRIA *et al.*, 2014). A pontuação polaridade (positivo e negativo) é um valor real (1,-1), semelhante a valores de polaridade fornecidos por outros métodos. O vector *Sentic* é composto por dezenas de emoções

⁴ *sentic computing*: A análise do texto não é baseada em modelos de aprendizagem estatísticos, mas sim em ferramentas de raciocínio de senso comum (CAMBRIA *et al.*, 2009) e ontologias específicas de domínio (CAMBRIA *et al.*, 2010).

agrupadas em quatro dimensões: simpatia, atenção, sensibilidade e aptidão. Estas dimensões baseiam-se no modelo *Hourglass emotions* (CAMBRIA et al., 2012). A versão 3.0 do *SenticNet*, está disponível em <http://sentic.net/>.

Destaca-se a utilização do *SenticNet* nos trabalhos (PORIA; GELBUKH; HUSSAIN, 2013; GEZICI et al., 2013).

EmoSenticNet. É um dicionário criado por (PORIA et al., 2014) que atribui seis conjuntos de emoções do *WordNetAffect* (raiva, alegria, repugnância, tristeza, surpresa, medo) para os conceitos de senso-comum do *SenticNet*, fornecendo a polaridade dos conceitos a cada conjunto do *WordNetAffect*, conforme apresentado na Tabela 2.11. O léxico pode ser visto como uma extensão do *WordNetAffect* para um vocabulário maior para a língua inglesa (EMOSENTICNET, 2015). O *EmoSenticNet* contém mais de 13.000 conceitos, incluindo aqueles já existentes na lista do *WordNetAffect*. O léxico *EmoSenticNet* está disponível em (EMOSENTICNET, 2015).

Tabela 2.11: Amostra do léxico *EmoSenticNet* (EMOSENTICNET, 2015).

Conceito	Raiva	Desgosto	Alegria	Triste	Surpresa	Medo
peace	0	0	1	0	0	0
indifference	0	1	0	0	0	0
impatience	1	0	0	0	0	0
flurry	1	0	0	1	0	0
where	0	1	0	0	1	0
emergency	0	0	0	0	0	1

EmoSenticSpace. Com a necessidade de obter bases de conhecimento com informações semânticas, conceituais e afetivas, os autores em (PORIA et al., 2014) propuseram a criação do *EmoSenticSpace*, que torna-se uma extensão do *WordNetAffect* e *SenticNet*, através do fornecimento de ambas as emoções contidas nos léxicos e suas polaridades.

Para construir o *EmoSenticSpace*, os autores utilizaram a combinação de dois conjuntos, sendo eles: *ConceptNet* e *EmoSenticNet*. O primeiro conjunto refere-se a um projeto de representação do conhecimento, por meio de um grande grafo semântico que descreve o conhecimento humano em geral e como esse conhecimento é expresso em linguagem natural (SPEER; HAVASI, 2013). Os nós do grafo representam conceitos e as

arestas relações. Por exemplo, os conceitos *colher* e *comer* na relação ‘*usado para*’ resulta em (*colher – usado para – comer*) e os conceitos *livro* e *papel* na relação ‘*feito de*’ resulta em (*livro – feito de – papel*). Para mais informações sobre o *ConceptNet*, sugere-se a leitura de (SPEER; HAVASI, 2013). O dicionário *EmoSenticNet* é descrito nos parágrafos anteriores, e representado nesse trabalho como um grafo direcionado, de forma semelhante ao *ConceptNet*.

Para o método de combinação dos conjuntos, foram transformados os grafos em matrizes e aplicada a técnica de *Blending* para a inferência entre as duas matrizes. Essa técnica permite que múltiplas matrizes possam ser combinadas em uma única matriz, baseando-se na sobreposição entre essas matrizes. Como resultado uma nova matriz foi criada, contendo grande parte das informações compartilhadas pelas duas matrizes originais, ainda, para descartar os componentes que representam variações relativamente pequenas nos dados, foi submetido o método *Truncated Singular Value Decomposition (TSVD)* (PORIA; GELBUKH; AGARWAL et al., 2013).

SentiLex-PT. É um léxico de sentimento desenvolvido por (SILVA et al., 2012) pertencente ao grupo de pesquisa *Data Management and Information Retrieval*⁵ (DMIR) de Lisboa, com o intuito de conceber um léxico para a análise de sentimento e opinião sobre entidades humanas em textos redigidos em português de Portugal. Trata-se de um recurso pioneiro para esta língua, sendo atualmente constituído por 7.014 lemas e 82.347 formas flexionadas. Os atributos de sentimentos descritos em cada entrada do léxico são: (i) o alvo da manifestação de sentimento, (ii) polaridade do sentimento e (iii) o método de atribuição de polaridade (CARVALHO; SILVA, 2015).

Segundo os autores, a informação de polaridade associada às entradas foi, na maioria dos casos, manualmente atribuída. Certas entradas adjetivais foram, contudo, automaticamente classificadas por uma ferramenta (denominada *Judgment Analysis Lexicon Classifier - JALC*) desenvolvida para este fim. As formas flexionadas dos verbos e das expressões idiomáticas, bem como os respetivos atributos morfológicos foram extraídos semi-automaticamente do LABEL-LEX, um léxico de palavras simples desenvolvido pela equipe do LABEL⁶ (RANCHHOD et al. 1999) para o português de Portugal.

⁵ http://dmir.inesc-id.pt/project/Main_Page

⁶ http://label.ist.utl.pt/pt/labellex_pt.php

O léxico é disponibilizado em arquivo de texto CSV⁷, deixando a cargo da aplicação “cliente” processar as informações do arquivo e representar os dados em um formato funcional próprio (MOURÃO; SAIAS, 2013). A Figura 2.4 ilustra uma pequena amostra do arquivo.

```
aberração.PoS=N;TG=HUM:N0;POL:N0=-1;ANOT=MAN
bonito.PoS=Adj;TG=HUM:N0;POL:N0=1;ANOT=MAN
castigado.PoS=Adj;TG=HUM:N0;POL:N0=-1;ANOT=JALC
estimado.PoS=Adj;TG=HUM:N0;POL:N0=1;ANOT=JALC;REV=AMB
enganar.PoS=V;TG=HUM:N0:N1;POL:N0=-1;POL:N1=0;ANOT=MAN
engolir em seco.PoS=IDIOM;TG=HUM:N0;POL:N0=-1;ANOT=MAN
```

Figura 2.4: Amostragem do arquivo CSV do SentiLex-PT (SENTILEX-PT, 2015).

Em cada linha do arquivo é dada as informações sobre: (i) Lema (convencionalmente a forma masculina do singular para os adjetivos, a forma singular para os nomes que flexionam em número e a forma infinitiva para os verbos e expressões idiomáticas), (ii) categoria gramatical, como: ADJ(etivo), N(ome), V(erbo) e IDIOM(a) e (iii) atributos de sentimento: polaridade(POL), a qual pode ser positiva (1), negativa (-1) ou neutra (0); alvo da polaridade (TG), o qual corresponde a um nome de tipo humano (HUM), com função de sujeito (N0) e/ou complemento (N1); classificação de polaridade (ANOT), a qual pode ter sido manualmente (MAN) ou automaticamente atribuída, pela ferramenta (JALC), desenvolvida pela equipe do projeto. Algumas entradas incluem um código adicional (REV), o qual se refere a observações específicas do anotador. Neste momento, é possível encontrar as seguintes anotações: (i) REV=AMB, associada a entradas cujo predicador é ambíguo com outra expressão que apresenta uma polaridade diferente, e (ii) REV:POL, associada a entradas cujo valor de polaridade inicialmente atribuído na versão inicial do SentiLex-PT foi revisto (SENTILEX-PT, 2015).

OpLexicon. Léxico de sentimento criado por (SOUZA et al., 2011) para o português brasileiro, constituído de 31.144 palavras classificadas segundo sua categoria morfológica e anotado com polaridades positivas, negativas e neutras (1, -1 e 0 respectivamente). Ele foi construído por fontes de informação, tais como: *corpus* (textos jornalísticos e resenhas de filmes), thesaurus (TEP thesaurus (MAZIERO et al., 2008)) e a tradução do léxico *English*

⁷ CSV: *comma-separated values*, é um formato onde há um conjunto de valores em cada linha, separados por vírgula ou outro separador textual como “:” ou “.”.

Opinion Lexicon, criado por (Hu; Liu 2004). A partir do processamento dos dados, três léxicos de opinião diferentes foram unidos para criar um grande léxico para o português brasileiro. A Tabela 2.12 apresenta a estrutura do léxico *OpLexicon*, onde a coluna conteúdo é composta por: palavra + adjetivo ou verbo + polaridade (1,-1 ou 0), separado por vírgula.

O *OpLexicon* está presente nos trabalhos de (SOUZA; VIEIRA, 2012; FREITAS; VIEIRA, 2013).

Os trabalhos de Mairesse (MAIRESSE, 2007) e Celli (CELLI, 2012) mostraram que as análises baseando em características linguísticas e palavras emotivas, carregam grande potencial para identificar a personalidade associada ao texto. Essa abordagem influenciou demais pesquisadores que propuseram a adição de palavras emotivas para a tarefa da identificação de personalidade. Em (MOHAMMAD; KIRITCHENKO, 2012), os autores compararam os resultados obtidos por Mairesse e acrescentaram no experimento os léxicos afetivos de emoções ligadas às palavras com *hashtag* (#love, #annoyed, #pity) no *Twitter*, e palavras anotadas manualmente para a detecção de personalidade, descobrindo assim um melhoramento essencial na precisão do sistema de RPT. Ainda, (PORIA; GELBUKH; AGARWAL et al., 2013) propôs uma nova arquitetura para a tarefa de reconhecimento de personalidade utilizando emoções com o uso do conhecimento de senso comum presentes nos léxicos *ConceptNet*, *EmoSenticNet*, *EmoSenticSpace*, junto com os léxicos LIWC e MRC.

Da mesma maneira, o presente trabalho usará léxicos afetivos com o intuito de melhorar a precisão do reconhecimento de personalidade, no entanto, se faz diferença na utilização de léxicos afetivos em língua portuguesa.

Tabela 2.12: estrutura do léxico *OpLexicon* (OPLEXICON, 2015).

Nº Linha	Conteúdo
8	abafado,adj,-1
4479	batizar,vb,1
5955	cego,adj,-1
6311	cidadão,adj,0
14212	estampar,vb,1
16051	fosco,adj,-1
20018	ler,vb,0
22107	música,adj,1

As abordagens baseadas em léxicos são as habitualmente utilizadas para resolver o problema de inferência de personalidade em textos. Todavia, há autores que utilizam-se de outros métodos, além dos léxicos supracitados. A próxima seção apresenta as demais abordagens, ainda que menos usuais, aplicados na literatura.

2.2.3. Outras abordagens

Os estudos apresentados até o momento adotam abordagens lexicais, ou seja, eles se baseiam em palavras individuais pré-rotuladas, em sua maioria por seres humanos, contendo significado afetivo e dados linguísticos. Embora os léxicos sejam amplamente utilizados, há pesquisadores que se utilizam de outras abordagens para a inferência da personalidade por meio de texto.

Os autores (OBERLANDER; NOWSON, 2006) trabalharam na classificação automática de personalidade a fim de detectar quatro dos cinco grandes traços (Extroversão, Neuroticismo, Socialização e Realização) em um *corpus* de *blogs* pessoais. Eles exploraram a técnica de extração de característica do texto chamado de n-gram (CAVNAR; TRENKLE, 1994), na etapa de pré-processamento da mineração de texto.

Esse recurso é a representação de uma cadeia de palavras de tamanho n , isto é, o número de termos que irá compor a cadeia do n-gram, podendo ser: unigram, contendo apenas 1 termo; bigram, contendo 2 termos; trigram, contendo 3 termos; e outros. Por exemplo, a frase “ser ou não ser” e considerando $n = 2$ (bigram), terá como resultado uma cadeia “ser ou”, “ou não” e “não ser”. O objetivo principal da aplicação dessa técnica é a obtenção de sequências de palavras que formam termos únicos e prever o próximo item em uma sequência. Nesse contexto, os termos “text” e “mining” terão maior relevância juntos do que os termos analisados separadamente (SILVEIRA et al., 2011). Além de palavras, os itens do n-gram podem ser formados de fonemas ou sílabas, de acordo com a aplicação. Ainda, cabe ressaltar os estudos (NOWSON; OBERLANDER, 2007) e (LUYCKX; DAELEMANS, 2008), que também exploraram n-gram para a tarefa de reconhecimento da personalidade por meio de texto.

Em (QUERCIA et al., 2011), os autores inferiram a personalidade por meio de texto dos usuários da rede social *Twitter*, explorando as características de relacionamento dos usuários na rede, como por exemplo a quantidade de seguidores, perfis que o usuário segue e listas (categorias de Twitters de temas específicos). Com essas características foi possível

identificar três tipos de usuários: ouvintes (aqueles que seguem muitos usuários), populares (aqueles que são seguidos por muitos), e leitores (aqueles que são muitas vezes mencionados em listas de leitura dos outros usuários). Dessa maneira, os autores utilizaram-se de algoritmos de aprendizagem de máquina para criar um modelo de inferência de personalidade a partir de 335 usuários da rede social *Twitter*.

Outros autores também analisaram características dos perfis dos usuários das redes sociais. Em (BACHRACH et al., 2012), foi criado um modelo por meio de algoritmo de aprendizagem de máquina, utilizando-se características da rede social *Facebook*, como: o tamanho da rede de relacionamento do usuário (amigos), a postagem de fotos, a confirmação em eventos, as marcações em fotos e outros. Segundo os autores, essas características demonstraram uma correlação fraca com os traços de personalidade, pois possuem inúmeras limitações. Isto porque, as características que foram analisadas apresentam um alto nível de agrupamento, por exemplo, foram examinadas a quantidade de “likes” ao invés de quais objetos o usuário gostou, ou ainda, a quantidade de mensagens ao invés das palavras utilizadas nas mensagens.

De outra forma, os pesquisadores (ALAM; STEPANOV; RICCARDI, 2013) utilizaram-se da abordagem chamada TF-IDF (*Term Frequency - Inverse Document Frequency*) para extrair características dos textos e criar um modelo de inferência de personalidade. Esse método leva em consideração a frequência de um termo em um documento, ou seja, quanto mais frequente um termo ocorre em um documento, maior representatividade esse termo é para o conteúdo, ou ainda, quanto mais documentos o termo ocorre, menos distinto ele é para o conteúdo (SALTON; BUCKLEY, 1988). Esse processo auxilia a distinguir o fato da ocorrência de algumas palavras serem geralmente mais comuns que outras. A atribuição do valor TF-IDF para os termos é dado pela Equação 2.1, que primeiramente calcula a medida *Term Frequency* – TF – o qual representa o número de vezes que o termos t_j ocorre no documento d_i e posteriormente para atribuir peso aos termos é calculado a medida *Inverse Document Frequency* – IDF – que varia inversamente ao número de documentos x que contem o termo t_j em um conjunto de documentos N .

$$TF - IDF(t_j, d_i) = tf(t_j, d_i) \times \log \frac{N}{x} \quad (2.1)$$

Ainda, neste estudo, os autores representaram o documento de forma vetorial utilizando o método *bag-of-words*, também chamado de Frequência de Características (*Feature Frequency*) por (PANG; LEE; VAITHYANATHAN, 2002). Esse método é a maneira mais simples de representar os documentos em um processo de classificação de textos, pois cria uma lista dos termos que aparecem no conjunto total de documentos e representa cada documento como um vetor composto pelo número de ocorrências de cada termo. Por exemplo, se o termo "personalidade" está presente três vezes em um determinado documento, o peso associado a esse termo no vetor será valor 3.

Outros pesquisadores também exploraram o método TF-IDF para a inferência de personalidade por meio de texto, por exemplo, os estudos de (PENG et al., 2015), (IACOBELLI et al., 2011) e (ARROJU; HASSAN; FARNADI, 2015). Desta maneira, este estudo também se utilizará do método TF-IDF para selecionar os termos com base em sua frequência e submetê-los aos léxicos afetivos, apresentado com mais informações no Capítulo 4.

Com base no levantamento bibliográfico apresentado ao longo da Seção 2.2, a Tabela 2.13 apresenta um resumo do estado da arte considerando algumas das principais e mais atuais abordagens, mencionadas a cima, para a inferência de personalidade por meio de textos.

Tabela 2.13: Resumo do estado da arte das abordagens de identificação de personalidade por meio de texto.

Trabalhos	Corpus	Extração de Características	Algoritmos ⁸	Traços				
				Ext.	Neu.	Soc.	Rea.	Abe.
Oberlander e Nowson (2006)	Blogs	N-gram	NB e SMO	Sim	Sim	Sim	Sim	Não
Mairesse et al. (2007)	Documentos	LIWC e MRC	LR, M5, REP e SMO	Sim	Sim	Sim	Sim	Sim
Rigby e Hassan (2007)	Emails	LIWC	C4.5	Sim	Sim	Sim	Sim	Sim
Gill et al. (2009)	Blogs	LIWC	LR	Sim	Sim	Sim	Sim	Sim
Yarkoni (2010)	Blogs	LIWC	Correlação	Sim	Sim	Sim	Sim	Sim
Iacobelli et al. (2011)	Blogs	LIWC, Bigram e TF-IDF	SVM, SMO e NB	Sim	Sim	Sim	Sim	Sim
Roshchina et al. (2011)	TripAdvisor	LIWC e MRC	LR, M5 e SVM	Sim	Sim	Sim	Sim	Sim
Quercia et al. (2011)	Twitter	Dados da Rede Social	M5	Sim	Sim	Sim	Sim	Sim
Golbeck et al. (2011)	Twitter	LIWC e MRC	ZeroR e Gaussian Processes	Sim	Sim	Sim	Sim	Sim

⁸ Nearest neighbour (KNN), Naive Bayes (NB), Linear Regression (LR), Máquina de Vetores de Suporte (SVM), Sequential Minimal Optimization (SMO), REPTree (REP) e Stochastic Gradient Descent (SGD).

Trabalhos	Corpus	Extração de Características	Algoritmos	Traços				
				Ext.	Neu.	Soc.	Rea.	Abe.
Bachrach et al. (2012)	Facebook	Dados da Rede Social	LR e SVM	Sim	Sim	Sim	Sim	Sim
Bai et al. (2012)	Rede Social Chinesa (RenRen)	Dados da Rede Social	NB, SVM e C4.5	Sim	Sim	Sim	Sim	Sim
Sumner et al. (2012)	Twitter	LIWC e Dados da Rede Social	SVM, NB, J48 e RandomForest	Narcisismo, Maquiavélico e Psicopata				
Markovikj et al. (2013)	Facebook	Dados da Rede Social	SMO	Sim	Sim	Sim	Sim	Sim
Zuo et al. (2013)	Documentos	Stanford Parser e HowNetKnowledge	KNN	Sim	Sim	Sim	Sim	Sim
Poria et al. (2013)	Documentos	ConceptNet, EmoSenticNet, EmoSenticSpace, LIWC e MRC	SMO	Sim	Sim	Sim	Sim	Sim
Tomlinson; Hinote e Bracewell (2013)	Facebook	WordNet	Correlação	Não	Não	Não	Sim	Não
Alam, Stepanov e Riccardi, (2013)	Facebook	TF-IDF	SMO	Sim	Sim	Sim	Sim	Sim
Farnadi et al. (2013)	Facebook	LIWC e Dados da Rede Social	SVM, KNN e NB	Sim	Sim	Sim	Sim	Sim
Poria et al. (2014)	Documentos	EmoSenticSpace	SVM	Sim	Sim	Sim	Sim	Sim
Verhoeven; Company e Daelemans (2014)	Documentos	N-gram e LIWC	SVM	Sim	Sim	Sim	Sim	Sim
Arroju, Hassan e Farnadi (2015)	Twitter	LIWC, TF-IDF	SGD	Sim	Sim	Sim	Sim	Sim
Peng et al. (2015)	Facebook	Unigram e TF-IDF	SVM	Sim	Sim	Sim	Sim	Sim
Litvinova, Seredin e Litvinova (2015)	Documentos	Parts-of-speech Bigram	Correlação	Sim	Sim	Sim	Sim	Sim

Após serem observados os principais estudos da área envolvendo léxicos e outras abordagens, cumpre mencionar que alguns pesquisadores concentram-se em expandir a aplicação de reconhecimento de personalidade em textos de outros idiomas diferentes do inglês, como por exemplo: Italiano (CELLI, 2012), Grego moderno (KERMANIDIS, 2012), Mandarim (BAI et al., 2012), Espanhol (ARROJU; HASSAN; FARNADI, 2015) e Russo (LITVINOVA; SEREDIN; LITVINOVA, 2015).

Todavia, para a língua portuguesa do Brasil, foram constatados poucos indícios de estudos computacionais com o objetivo de inferir a personalidade humana, a saber, (NUNES; TELES; DE SOUZA, 2013; LIMA; CASTRO, 2013). Ambos os estudos têm como objetivo a

inferência da personalidade utilizando textos enviados pela rede social *Twitter*, seus métodos de detecção são baseados em ferramentas de *text-mining*. Na próxima seção serão exploradas as etapas de cada estudo.

2.3. Inferência de Personalidade por meio de textos para o português do Brasil

Devido ao seu potencial de aplicabilidade na língua inglesa, tem-se observado o início das pesquisas para o português brasileiro, com o intuito de expandir o RPT a outras línguas. Nesse sentido, os estudos realizados por (NUNES; TELES; DE SOUZA, 2013) e (LIMA; CASTRO, 2013) são as únicas evidências computacionais com o objetivo de detectar a personalidade humana em textos para o português do Brasil.

O trabalho de Nunes descreve um experimento baseado nas publicações de textos públicos da rede social *Twitter*, sendo que possui o escopo de buscar indícios que possibilitem comprovar correlações entre: ferramentas tradicionais de mensuração de personalidade, explicitamente por meio de inventários, e uma ferramenta baseada em *text-mining*, de forma implícita.

O referido autor fragmentou o experimento em duas partes: a primeira consistiu em capturar dos participantes (vinte e oito participantes) os textos/*tweets* públicos produzidos no *Twitter* e, em paralelo, construiu uma ferramenta que possibilitou a mineração desses textos inferindo traços de personalidade dos participantes relacionados ao modelo *BigFive*; e a segunda parte consistiu em coletar os traços de personalidade desses mesmos participantes usando dois inventários explícitos, tradicionais e reputados para a inferência da personalidade (a versão dos inventários utilizada foi em língua portuguesa).

Na etapa da mineração de texto, o *dataset* foi composto por amostras de *posts* dos vinte e oito participantes, coletadas de forma explícita pelos pesquisadores. Posteriormente, na fase do pré-processamento, os textos foram transformados em vetores de palavras e suas frequências foram sumarizadas utilizando técnicas de redução de dimensionalidade, como por exemplo, o *stemming*⁹.

Para auxiliar na sumarização dos termos durante o pré-processamento, um dicionário foi utilizado para uma rápida classificação e sumarização das palavras. O dicionário utilizado pelos autores foi o TeP 2.0 - *Thesaurus* Eletrônico para o Português do Brasil - (MAZIERO et

⁹ *Stemming* é um processo que consiste em reduzir a palavra ao seu radical.

al., 2008) que propõe uma análise sintática das sentenças. Dessa forma, o dicionário foi adaptado para separar as categorias das palavras em diversos tipos de análises, sendo que essa separação, segundo os autores, permitiu facilmente adicionar ou remover categorias de palavras, tornando, assim, o dicionário mais flexível. Essa adaptação foi baseada no modelo LIWC da língua inglesa (PENNEBAKER et al., 2001). A Tabela 2.14 mostra um exemplo da catalogação de alguns termos do dicionário.

Tabela 2.14: Exemplo de termos do dicionário produzido pelos autores (NUNES; TELES; DE SOUZA, 2013).

Dimensões	Exemplo de Palavras
Pronomes	Eu, nosso, eles, vocês
1ª pessoa singular	Eu, meu, mim
Negação	Não, nunca, nenhum
Preposições	Em, para, de
Números	Um, trinta, milhão
Processo afetivo ou emocional	Feliz, feio, amargo
Emoções positivas	Feliz, bastante, bom
Sentimentos positivos	Feliz, alegria, amor

Para mapear a personalidade dos participantes ao modelo *BigFive*, os autores utilizaram padrões apontados por Mairesse (MAIRESSE et al., 2007), que propõem uma expressão de mapeamento para o fator de Extroversão (E) como sendo a soma das frequências dos termos da seguinte forma:

$$E = (\text{substantivo} + \text{adjetivo} + \text{preposição} + \text{artigo} - \text{pronomes} - \text{verbos} - \text{advérbios} - \text{interjeição} + 100) / 2$$

Seguindo uma lógica semelhante, os autores propuseram as seguintes expressões para cada um dos fatores: Socialização (A), Realização (C), Neurotismo (N) e Abertura (O):

- $A = (\text{palavras positivas} - \text{palavras negativas} - \text{artigos} + 100) / 2$.
- $C = (\text{palavras positivas} - \text{negações} - \text{palavras negativas} + 100) / 2$.
- $N = (\text{palavras negativa} - \text{palavras positivas} + \text{primeira pessoa} - \text{terceira pessoa} + 100) / 2$.
- $O = (\text{palavras longas} - \text{palavras curtas} / 2 + \text{terceira pessoa} - \text{primeira pessoa} + 100) / 2$.

Nota-se que para a elaboração das expressões os autores ocultaram a maneira com a qual foram concebidas as sequências de termos e as escolhas destes.

A segunda etapa do experimento foi a coleta explícita da personalidade dos participantes por meio do preenchimento dos inventários TIPI (*Ten-Item Personality Inventory*) e NEO-IPI (*Neo-International Personality Item Pool*), disponíveis na língua portuguesa. Posteriormente, esses dados foram comparados ao processo de inferência da personalidade realizada pela mineração de texto.

Os resultados obtidos mostram uma correlação fraca entre os padrões de personalidade dos participantes obtidos via *Text-Mining* em relação aos padrões obtidos via inventários (NEO-IPIP e TIPI). Observa-se na Tabela 2.15 a correlação dos resultados do *Text-Mining* com os resultados dos inventários NEO-IPIP.

Tabela 2.15: Correlação entre NEO-IPIP & Text-Mining, adaptado de (NUNES; TELES; DE SOUZA, 2013).

<i>Text-Mining</i> / NEO-IPIP	Extroversão	Socialização	Realização	Neuroticismo	Abertura à experiência
Extroversão	0,1267	0,2441	0,1464	0,0144	0,1055
Socialização	-0,0282	-0,1611	0,0706	-0,1274	-0,0382
Realização	-0,2260	-0,4360	-0,1417	0,4122	0,2076
Neuroticismo	0,2056	0,4176	0,1723	-0,0314	0,1982
Abertura à experiência	-0,0734	0,0493	-0,0282	-0,1445	-0,0271

Os dados da correlação entre *Text-Mining* e os inventários TIPI, são mostrados na Tabela 2.16.

Tabela 2.16: Correlação entre TIPI & Text-Mining, adaptado de (NUNES; TELES; DE SOUZA, 2013).

<i>Text-Mining</i> / TIPI	Extroversão	Socialização	Realização	Neuroticismo	Abertura à experiência
Extroversão	0,2253	0,3535	0,1389	0,2593	0,2705
Socialização	0,1545	-0,2559	-0,0203	-0,0555	-0,1806
Realização	-0,2431	-0,0613	-0,2262	0,0384	-0,3753
Neuroticismo	0,4161	0,2849	0,2758	0,1722	0,2764
Abertura à experiência	0,0842	0,2132	-0,0397	-0,2206	-0,0627

Como pode ser observado, os dados apresentados demonstram que existe uma correlação fraca entre o *Text-Mining* e TIPI no fator Extroversão e Neuroticismo; e a correlação fraca entre o *Text-Mining* e o NEO-IPIP no fator Extroversão.

Os resultados obtidos por (NUNES; TELES; DE SOUZA, 2013) foram promissores, pois mostraram algum tipo de correlação, ainda que não satisfatória. Entretanto, a existência de uma correlação, mesmo que insatisfatória, estimula pesquisadores a aperfeiçoar os métodos para atingir melhores resultados.

Cabe, ainda, mencionar o estudo de (LIMA; CASTRO, 2013) que descreve também um experimento baseado nas publicações da rede social *Twitter*, a fim de mensurar a personalidade em texto na língua portuguesa. Ao contrário de Nunes, a análise de Lima e Castro se aplica em prever traços de personalidade dominante a um grupo de *tweets*¹⁰ sobre um determinado assunto, ao invés de exclusivamente prever traços de um único usuário. Para isso, um modelo de classificação semi-supervisionado (empregando estratégia multi-rótulo) foi desenvolvido e avaliado para prever a personalidade de acordo com o *BigFive*.

O estudo de Lima e Castro está dividido em três etapas: a primeira etapa é a fase do pré-processamento, em que são extraídas características dos textos. A lista de características utilizadas é dada a seguir, pelo número médio de:

- Seguidores;
- Seguintes;
- Palavras;
- Perguntas;
- Pontos de exclamação;
- Palavras com mais de seis letras;
- Palavras positivas;
- Palavras negativas;
- Artigos;
- Links na mensagem;
- Dois pontos (:).

A partir da coleta dessas informações, um novo conjunto de dados é criado, denominado *meta-base*. A próxima etapa é a fase de transformação, responsável em aplicar o conjunto *meta-base* a uma abordagem multi-rótulo, utilizando o algoritmo de classificação

¹⁰ *Tweet* é a mensagem postada no Twitter e possui um comprimento máximo de 140 caracteres.

Naïve Bayes, convertendo um problema de classe única, ou seja, cada dimensão do modelo *BigFive* é separado em um classificador binário (tem ou não o traço de personalidade). O sistema foi chamado de *Bayesian Prediction Personality* (ppBayes).

Por último, na terceira etapa, o classificador recebe os cinco conjuntos de treinamento e o conjunto *meta-base* (formado na primeira etapa). Com um método de classificação semi-supervisionado, o sistema classifica as amostras do *meta-base* e acrescenta o resultado para o conjunto de treinamento. A Figura 2.5 ilustra este procedimento em que o aprendizado utiliza dados marcados e não marcados no processo de previsão.

Os autores coletaram para o experimento *tweets* sobre opiniões ligadas a programas de televisão de estações brasileiras, com o propósito de compreender o comportamento coletivo sobre tais programas nos meios de comunicação social. Cerca de 4.621 *tweets* foram coletados e divididos em grupos de acordo com o nome do programa.

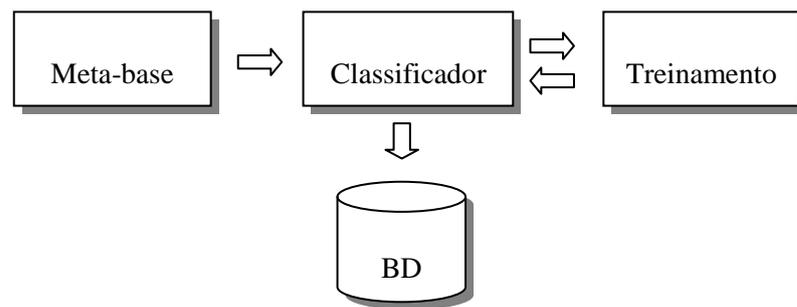


Figura 2.5: Processo de classificação semi-supervisionado utilizado por (LIMA; CASTRO, 2013).

A Tabela 2.17 ilustra o pré-processamento dos grupos, formando o conjunto *meta-base* com 30 objetos e 11 atributos. Em seguida, a Tabela 2.18 apresenta classificações estabelecidas pelo ppBayes.

Os pesquisadores afirmam que este tema foi abordado devido à demanda por sistemas de TV Sociais, e que os resultados apresentados mostram que o ppBayes é capaz de estimar, com uma acurácia de 84,15% e precisão de 87,37%, traços de personalidade em textos. O grande desafio de pesquisa para a inferência de personalidade em textos para o português é estabelecer novas relações entre os atributos e características de personalidade.

Tabela 2.17: Pré-processamento dos grupos, formando o conjunto meta-base com 30 objetos e 11 atributos, adaptado de (LIMA; CASTRO, 2013).

	TV Fama	Programa da Tarde	Encontro Fátima	Salve Jorge	Sangue Bom
Seguintes	989,85	160,22	665,65	10,52	58,46
Seguidores	2087,75	6726,88	97219,96	853,72	2548,29
Palavras	18,75	16,79	21,70	15,28	16,92
Perguntas	0,35	0,10	0,01	0,16	0,30
Exclamação	0,17	0,28	0,71	0,39	0,37
Palavras superiores a seis letras	5,18	6,59	8,44	4,97	6,49
Positivo	0,11	0,21	0,26	0,18	0,20
Negativo	0,09	0,03	0,00	0,07	0,14
Artigos	1,54	0,97	1,01	0,84	1,08
Links	0,18	0,53	0,20	0,13	0,30
Dois pontos	0,78	0,74	1,31	0,71	1,13

Tabela 2.18: Amostra de classificação do ppBates, adaptado de (LIMA; CASTRO, 2013).

Programa de TV	Traços de Personalidade
MTV sem Vergonha	Abertura à experiência
Encontro Fátima	Extroversão, Realização, Abertura à experiência
Sangue Bom	Extroversão

2.4. Considerações Finais

Neste capítulo foram apresentados os conceitos básicos de personalidade, bem como a abordagem dos traços e os principais modelos criados para classificar a personalidade de um indivíduo, sendo eles: o modelo dos 16 fatores, modelo dos três “Superfatores” e o modelo dos cinco fatores (*BigFive*).

Além disso, foram apresentados os estudos de inferência de personalidade em língua inglesa e língua portuguesa do Brasil, bem como os principais léxicos, com os quais pretende-se desenvolver um método automático para detecção de personalidade por meio de textos para o português do Brasil, diminuindo a lacuna encontrada no estado da arte para tal inferência.

Ressalta-se que a proposta não é trivial, devido a escassez de léxicos e métodos para a língua portuguesa, comparado com os avanços da inferência de personalidade para outras línguas, como por exemplo, o inglês.

Foram apresentadas ainda outras abordagens, além de léxicos, comumente utilizadas para a inferência de personalidade por meio de texto usada na literatura.

O capítulo seguinte apresenta o método proposto para o reconhecimento de personalidade a partir de textos para o português, tema deste estudo.

Capítulo 3

Um Modelo para Inferir a Personalidade a partir de Textos em Língua Portuguesa

No Capítulo 2 foram apresentadas as dificuldades dos métodos para a tarefa de inferência da personalidade a partir de textos redigidos em português do Brasil. Dentre elas, destaca-se o grande desafio de estabelecer relações entre os atributos e as características de personalidade (LIMA; CASTRO, 2013). Este estudo apresenta um modelo de inferência de personalidade utilizando léxicos, associados a um método de representação de termos com base em sua frequência (TF-IDF), que visa estabelecer relações às características de personalidade, pretendendo atingir melhores resultados do que os apresentados na literatura.

Neste capítulo é apresentado um método detalhado para a inferência da personalidade do indivíduo por meio de textos escritos em língua portuguesa. Nesse contexto, a Seção 3.1 apresenta uma visão geral do método proposto. Posteriormente, a Seção 3.2 apresenta a aplicação explícita do inventário de personalidade, cujo resultado será utilizado para a avaliação do método. Logo após, na Seção 3.3, é mostrada a etapa de coleta dos textos. Por fim, a Seção 3.4 apresenta o pré-processamento dos dados e os léxicos que serão utilizados na criação do modelo.

3.1. Visão Geral do Método

Apresenta-se um método para a inferência da personalidade do indivíduo por meio de textos escritos em mídias sociais em língua portuguesa. A seção apresenta uma visão geral do modelo proposto e uma sucinta descrição sobre cada etapa do projeto.

Inicialmente foram selecionados para o experimento voluntários que possuem assiduidade de publicações em redes sociais, e a eles foi aplicado um questionário de Psicologia que visa extrair seus traços de personalidade. O objetivo de obter preliminarmente os traços de personalidade dos participantes, por meio de um inventário robusto e respaldado na área da Psicologia, permitirá a avaliação do método proposto. Na Seção 3.2 é apresentada com mais detalhes a mensuração explícita de personalidade por meio de inventário.

Posteriormente, para a montagem da base textual, são capturadas as mensagens publicadas na rede social dos participantes, o que permitirá medir suas dimensões de personalidade por meio dos textos postados. A Seção 3.3 descreve em detalhes a coleta dos dados.

Em seguida, a fase do pré-processamento é responsável pela extração de características linguísticas e emocionais empregadas nos textos. Para a tarefa de extração de características linguísticas, será utilizado o léxico LIWC que possui adaptação para o português brasileiro e que contém uma gama textual de categorias, por exemplo: advérbios, palavras de negação, tempos verbais, pronomes, pontuações e outros. Nesta fase será criada uma base de dados para cada traço de personalidade do modelo *BigFive*, isto é, há cinco bases de dados com atributos (características linguísticas, afetivas e frequência de termos) e exemplos rotulados com referência a um traço de personalidade do *BigFive*, resultando, assim, na construção das bases de dados para: Extroversão, Neuroticismo, Socialização, Realização e Abertura à experiência.

Para a tarefa de extrair emoções empregadas em textos, serão utilizados os léxicos afetivos: *SentiStrength*, *AnewBr*, *SentiLex-PT* e *OpLexicon* com o intuito de aperfeiçoar a associação da personalidade em textos. Todos os léxicos afetivos citados foram construídos ou adaptados para a língua portuguesa por outros pesquisadores. Ainda, justifica-se a utilização de vários léxicos com o propósito de comparar quais possuem melhores resultados, podendo ser combinados entre si.

Além disso, com o intuito de extrair os termos que possuem maior representatividade no texto e associá-los aos léxicos, para aperfeiçoar o método de inferência da personalidade, será utilizado o método TF-IDF. Maiores detalhamentos sobre as métricas de utilização do LIWC, léxicos afetivos e TF-IDF, estão disponíveis na Seção 3.4.

Para melhor entendimento do modelo proposto, a Figura 3.1 ilustra as fases do projeto e como elas estão relacionadas. Nota-se que após a etapa do pré-processamento são

construídas as bases de treinamento para os algoritmos de aprendizagem de máquina, a partir de características textuais e informações emocionais. O atributo meta (classe) é rotulado com o resultado do questionário de personalidade.

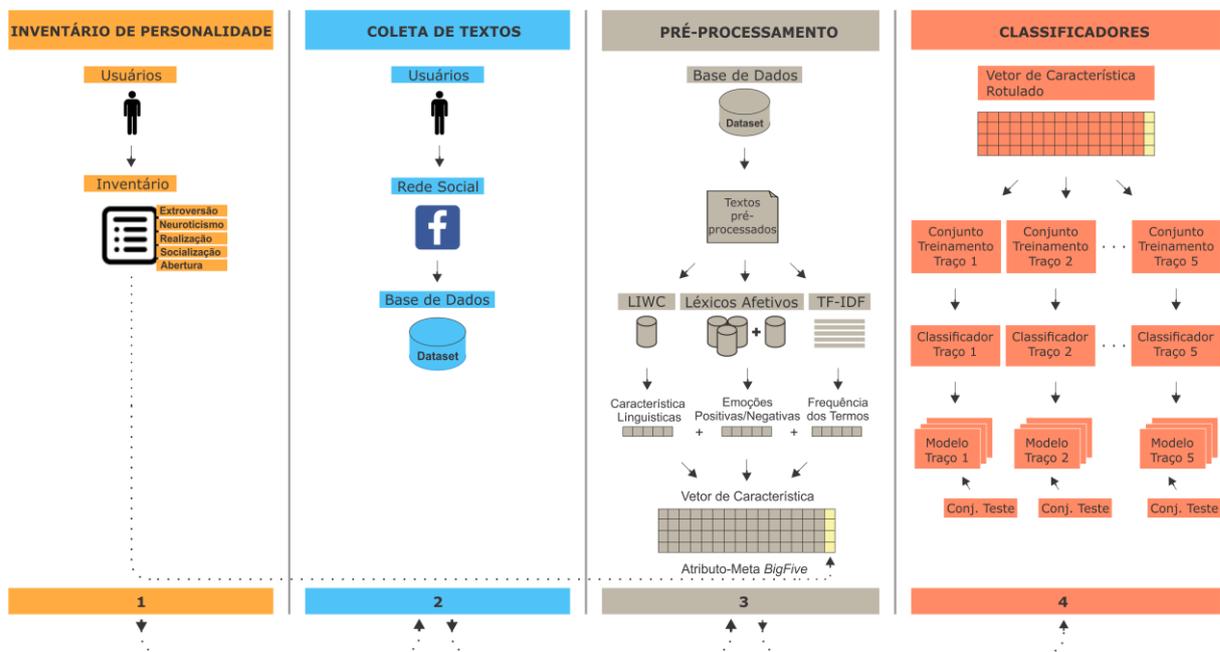


Figura 3.1: Visão geral do modelo para inferir a personalidade a partir de textos.

Observa-se que os métodos de inferência de personalidade a partir de textos em língua portuguesa, apresentados na literatura, utilizam limitadas características textuais, como observado no trabalho de (LIMA; CASTRO, 2013). A utilização do LIWC no estudo apresentado permite abranger um maior número de características textuais empregadas nos textos. Dessa forma, o uso de léxicos no modelo tornar-se um diferencial para esse estudo.

Além disso, observam-se outras limitações nos estudos existentes aplicados à língua portuguesa, como abstenções das emoções empregadas nos textos e avaliação do método a uma exígua quantidade de participantes, como por exemplo, o trabalho de (NUNES; TELES; DE SOUZA, 2013).

Referente ao processo de mineração dos dados, são utilizados algoritmos de regressão para geração de modelos de reconhecimento da personalidade, conforme as dimensões do *BigFive*. Os indutores de regressão serão utilizados pelo fato do atributo meta assumir valores contínuos. Posteriormente, serão confrontados entre si os resultados obtidos pelos algoritmos para a avaliação do método.

Cada etapa do modelo proposto será apresentada detalhadamente nas próximas seções.

3.2. Mensuração Explícita de Personalidade via Inventário

Na Seção 2.1 foram apresentados a abordagem de traços e os modelos de personalidade. Dentre esses modelos de personalidade destaca-se o *BigFive*, que descreve e classifica a personalidade humana em cinco fatores. Deste modo, utilizar-se-á o modelo *BigFive* para descrever a personalidade em torno dos cinco traços, sendo eles:

- Extroversão (Sociável *versus* Tímido);
- Neuroticismo (Calmo *versus* Neurótico);
- Socialização (Amigável *versus* Não cooperativo);
- Realização (Organizado *versus* Descuidado);
- Abertura à experiência (Intelectual *versus* Sem imaginação).

Conforme observado no segundo capítulo, o modelo *BigFive* é amplamente utilizado pelos pesquisadores na inferência de personalidade por meio de texto, já que, quando comparado com outro modelo, por exemplo, modelo dos três “superfatores” de Eysenck, ele descreve de forma clara e objetiva os traços de personalidade. Ainda, leva-se em consideração que o modelo possui grande aceitabilidade não apenas na área da computação afetiva, mas também, se relaciona com um vasto número de estudos no âmbito científico.

Desse modo, com a finalidade de extrair os traços de personalidade dos participantes, de maneira explícita, por meio de um inventário robusto e respaldado na área da Psicologia, o psicólogo¹¹ que cooperou com o estudo avaliou o questionário de personalidade NEO-IPIP, criado por Johnson (2014), contendo 120 questões.

A aplicação do questionário tem o intuito de auxiliar a avaliação do método proposto e rotular os exemplos do conjunto de treinamento, caracterizando uma aprendizagem supervisionada.

O inventário NEO-IPIP 120 foi destacado entre os demais presentes na literatura, pelo psicólogo, por abranger com precisão as representações dos traços de personalidade do modelo *BigFive* e possuir um número de questões viável para o preenchimento em ambiente computacional. Observa-se que, muitas vezes, a aplicação de um questionário com grande número de questões se torna inviável, pois requer do usuário o emprego de bastante tempo, o que pode levar à desistência da conclusão do teste, deixando o processo de avaliação do

¹¹ O profissional de Psicologia que acompanhou o projeto chama-se Ricardo Stegh Camati.

comportamento vulnerável. Assim, a curta versão do NEO-IPIP fornece uma alternativa para esse tipo de situação.

A utilização do inventário NEO-IPIP 120 também contou com a colaboração de Johnson, que autorizou a aplicação do mesmo, compartilhou métricas para computar cada questão do inventário e associou as questões aos traços do modelo *BigFive*. A versão original do inventário encontra-se em língua inglesa, contudo, o estudo de Nunes (2013), descrito na seção 2.3, realizou a tradução do inventário NEO-IPIP contendo 300 questões, e após a criação do inventário NEO-IPIP 120, a autora e sua equipe também realizaram a tradução, criando uma versão brasileira (língua portuguesa do Brasil) do inventário. Dessa forma, a partir da tradução de Nunes, foi possível criar uma interface computacional e impressa do inventário NEO-IPIP 120 para o português brasileiro, tendo como objetivo oferecer mais usabilidade ao questionário e acessibilidade às pessoas envolvidas nesta pesquisa.

Desse modo, o questionário NEO-IPIP 120 foi apresentado aos participantes dessa pesquisa na versão impressa e *online* para seu preenchimento. Foram mantidos o número de questões e o formato das respostas em ambas as versões, impressas e *online*, e asseguradas as características originais da versão criada por Johnson (2014).

Os voluntários do estudo responderam cada questão, selecionando uma das cinco opções de resposta que representam seu nível de concordância para cada pergunta, formatado em uma escala tipo Likert (LIKERT, 1932), conforme ilustrado na Figura 3.2. No Anexo 1, são apresentadas todas as questões do inventário NEO-IPIP 120, incluídas nos formulários *online* e impressos, contendo a relação entre as questões e os fatores do *BigFive*.

1. Preocupo-me com as coisas *

Discordo Totalmente	Discordo Parcialmente	Nem discordo nem concordo	Concordo Parcialmente	Concordo Totalmente
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Figura 3.2: Questão 1 do Teste NEO-IPIP 120. Adaptado de (JOHNSON, 2014).

Por fim, após o término do preenchimento total do questionário, os valores atribuídos a cada uma das indagações respondidas são utilizados para contabilizar o resultado. No cálculo, ao resultado é atribuído um valor entre 1 a 100 para cada um dos itens do *BigFive*. O valor 1 corresponde ao nível mais baixo do traço e o valor 100 ao nível mais alto do traço. Por

exemplo, um valor superior ou igual a 70 para o traço Extroversão, segundo o psicólogo que auxiliou na aplicação, identifica um indivíduo com um forte engajamento por emoções positivas e pela tendência em procurar estimulação e a companhia dos outros.

A seguir, na próxima subseção, será exposto o perfil dos participantes que preencheram o questionário de personalidade.

3.2.1. Participantes

Os participantes da pesquisa são estudantes universitários que possuem assiduidade nas mídias sociais e que responderam de forma voluntária o inventário NEO-IPIP 120, por meio do documento *online* ou impresso.

Antes do preenchimento do inventário, os estudantes anuíram com a participação no experimento (ver apêndice A), sendo que as informações relativas às respostas feitas no questionário serão mantidas sob sigilo e utilizadas exclusivamente para fins dessa pesquisa científica. Para a colaboração com o estudo, também foram solicitadas informações sobre a utilização das mídias sociais, tais como: nome do usuário, idade, sexo e a frequência de utilização das redes sociais, mais especificadamente o *Facebook*. As referidas informações foram utilizadas para coleta dos textos publicados nas redes sociais. Esse tema será discutido na próxima seção.

Com o preenchimento da autorização, obtiveram-se os seguintes dados dos participantes:

Tabela 3.1: Informações sobre os participantes do experimento.

Informações	Total
Total de participantes	575
Sexo masculino	351
Sexo feminino	224
Idade média	23

Todavia, apesar do grande número de integrantes que consentiram na participação do experimento, houve um filtro para a montagem do conjunto de dados. Foi analisado se os participantes que responderam o questionário possuem contas ativas nas redes sociais e extratos de dados textuais que permitam a extração e análise textual do conteúdo.

Desde modo, foram descartados do experimento os questionários: (i) incompletos em preenchimentos impressos, (ii) de usuários com perfil desativado no *Facebook* (iii) de usuários do *Facebook* não assíduos na rede (sem postagem), (iv) de usuários que não anuíram com o aplicativo para coleta dos dados e, (v) de usuários da rede social com exíguas publicações e que não permitiram a extração e análise textual do conteúdo, seguindo a metodologia utilizada em (PENG et al., 2015), com o fator de eliminação de usuários com menos de 1.000 palavras.

Após a validação desses itens, foi possível estabelecer quais participantes fariam parte da base de dados. Nesse contexto, a Tabela 3.2 apresenta o número final de integrantes que terão os dados coletados da rede social *Facebook*.

Tabela 3.2: Informações sobre os participantes finais do experimento.

Informações	Total
Total de participantes	256
Sexo masculino	115
Sexo feminino	141
Idade média	25

Ressalta-se que, exceto os questionários que contém respostas incompletas, todos os demais participantes do experimento receberam seu resultado de personalidade e uma breve explicação sobre cada traço do modelo *BigFive*.

A média do resultado de cada traço de personalidade dos participantes que compõe a base de dados é apresentada na Figura 3.3. A figura também apresenta a barra do desvio padrão para cada traço.

A análise dos resultados de cada traço, advindos do questionário, requer sólidos conhecimentos da área da Psicologia, dentre eles um aprofundado embasamento teórico do modelo *BigFive*. Por tal motivo, mesmo havendo acompanhamento de um profissional da Psicologia, a justificativa dos resultados apresentados não será explanada por não fazer parte do escopo da pesquisa. Tal tarefa poderá futuramente ser realizada com mais expertise por profissionais da área da Psicologia.

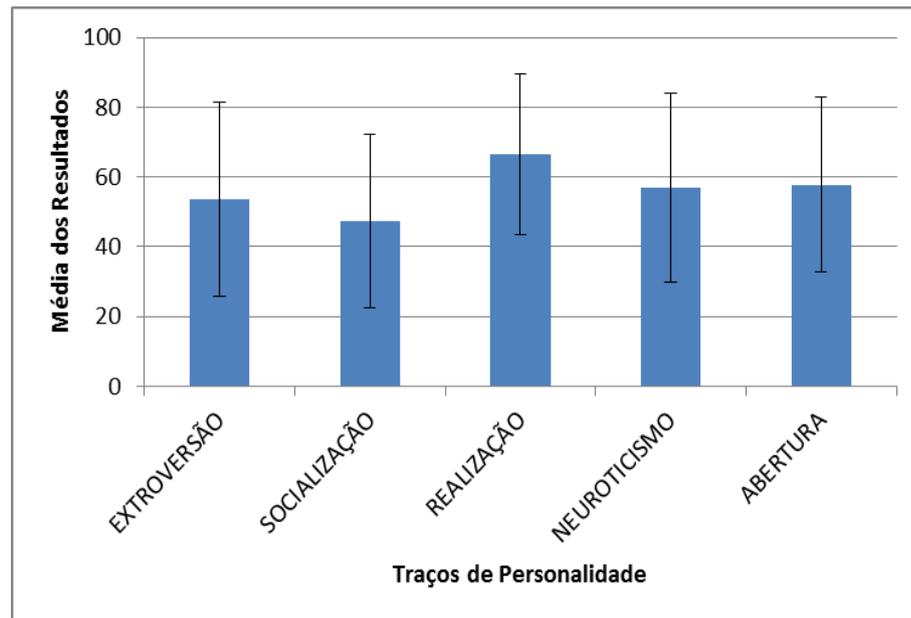


Figura 3.3: Média dos resultados com desvio padrão dos traços de personalidade via inventário.

Ressalta-se que o inventário aplicado aos participantes, bem como a pesquisa proposta neste trabalho, foram encaminhados para o Comitê de Ética em Pesquisa da PUCPR, que emitiu um Certificado de Apresentação para Apreciação Ética (CAAE), vide Anexo 2.

Após a inferência da personalidade por meio de inventário, foi realizada a coleta de textos publicados na rede social dos participantes selecionados. A próxima seção tratará as minúcias da montagem da base textual.

3.3. Base de Dados

É sabido que as opiniões geradas em redes sociais tornaram-se recentemente um recurso valioso para análise do comportamento dos usuários. Com o advento das mídias digitais, a comunicação e a informação se tornaram cada vez mais acessíveis e suas disseminações mais abrangentes, transformando a maneira do ser humano se relacionar entre sua comunidade. A rede social digital possui forte influência nesse aspecto, contribuindo na expansão da produção da informação e da comunicação, possibilitando aos usuários não apenas ter acesso à informação, mas também, gerar a informação.

Para o estudo de caso desta pesquisa, serão utilizadas publicações da rede social *Facebook*. Justifica-se a escolha de tal rede social em decorrência à vasta utilização da mesma pelo público brasileiro. Corrobora-se tal fato tendo em vista que, recentemente, um estudo de

consumo de mídia realizado pelo governo brasileiro classificou o *Facebook* como sendo a rede social mais utilizada no Brasil (SECOM, 2015).

A base textual dessa pesquisa faz uso do volume de dados gerados pela referida rede social. O processo tem como objetivo capturar as mensagens publicadas pelos participantes (os mesmos da seção 3.2), de maneira a permitir a mensuração de suas dimensões de personalidade por meio dos textos postados.

Dessa forma, a base de dados é composta por textos publicados independentemente da data em que o participante postou a mensagem, pois de acordo com a literatura da Psicologia, a personalidade de uma pessoa é normalmente estável, mas não imutável (BERGER, 2003). Assim, uma publicação realizada há três anos pode ser significativa na construção do modelo.

Posteriormente à etapa de seleção e criação da base de dados, o processo avança para a etapa de pré-processamento, responsável pela limpeza e representação dos dados. Essa etapa será detalhada na próxima seção.

3.4. Pré-processamento

Considerando a finalização da coleta dos textos e a construção da base de dados, o pré-processamento dos dados se faz necessário para formatar e organizar de maneira adequada o modelo para a inferência de personalidade. Essa etapa é responsável pela preparação do texto para a mineração dos dados, e também, realiza a extração de características linguísticas e emocionais empregadas nos textos, criando uma matriz de termos relevantes para a inferência.

Comumente a fase do pré-processamento em mineração de texto consiste em três etapas: análise léxica, eliminação de termos considerados irrelevantes e normalização morfológica dos termos. Segundo Riloff (1995), dependendo da aplicação de domínio, essas etapas podem variar sua ordem ou simplesmente não ocorrer.

Na análise léxica é feita uma adaptação do texto, eliminam-se os sinais de pontuação, isolam-se os termos (tokenização) e convertem-se letras maiúsculas para minúsculas, estágios que foram realizados nesse estudo. Por exemplo, considere a frase retirada do *Facebook* “ *O sucesso é ir, de fracasso em fracasso, sem perder o entusiasmo !. #ToNaLutaBrasil*”. A frase resultante da aplicação dessa etapa é “*o sucesso é ir de fracasso em fracasso sem perder o entusiasmo #tonalutabrasil*”. Observa-se que a frase encontra-se sem as vírgulas, o ponto de exclamação e o ponto final.

A remoção de termos irrelevantes ou *stopwords* visa a retirada das palavras de pouca importância para a representatividade do texto em um processo de mineração de texto. Essa retirada de termos que não possuem um valor semântico na frase, não agrega conteúdo ao assunto do texto em questão. Geralmente esses vocábulos se encontram na classe gramatical de verbos auxiliares, pontuações, artigos e preposições. Todavia, para esse estudo serão mantidas algumas classes gramaticais (como por exemplo: verbos, verbos auxiliares e artigos), devido ao léxico LIWC categorizar essas palavras, influenciando o vetor de característica e consequentemente o modelo proposto. Adiante serão exploradas mais informações sobre a utilização o LIWC no método.

Existem palavras consideradas comumente irrelevantes no processo de pré-processamento, como por exemplo, o advérbio “ainda”, mantido para a catalogação do LIWC, por tal motivo, faz-se necessária a elaboração de uma lista de *stopwords*. A seguir a Tabela 3.3 apresenta alguns exemplos de *stopwords* utilizados para essa etapa.

Tabela 3.3: Exemplos de *stopwords*.

<i>stopwords</i>
a, agora, alguém, algum, alguma, algumas, alguns, ampla, amplas, amplo, amplos, ante, antes, ao, aos, após, aquela, aquelas, aquele, aqueles, aquilo, as, até, através, cada, coisa

A lista de *stopwords* foi retirada do site da Linguateca¹², que contém as palavras mais comuns na língua portuguesa, independentemente da classe gramatical, e fornece diversos recursos para tal tarefa.

Cabe ressaltar que a coleta de dados publicados nas redes sociais necessita de tratamentos adicionais aos citados, pois é comum encontrar, por exemplo, mensagens repetitivas geradas por um mesmo autor, o que acarreta dados redundantes que podem ser computados. Visando reduzir o esforço dos processos subsequentes, todas as repetições de conteúdo serão eliminadas do processo. Também, palavras que utilizam do símbolo *hashtags* (#) serão segregadas, ainda que existam estudos de extração de sentimentos para a língua inglesa que fazem uso somente das palavras precedidas de *hashtags*. Para mais detalhes sugere-se a leitura de (DAVIDOV et al., 2010; MOHAMMAD; KIRITCHENKO, 2012).

¹² <http://www.linguateca.pt/>

Outro desafio encontrado diz respeito ao vocábulo utilizado pelos usuários nas redes sociais, em sua maioria em caráter informal. Tal informalidade permite que abreviações e variações na escrita de uma palavra representem um mesmo significado, por exemplo, as expressões “vc”, “amg” e “vdd” correspondem aos termos “você”, “amigo” e “verdade” respectivamente, o que pode dificultar o tratamento dos textos. Também, outro fator importante são as diversidades linguísticas presentes no idioma português, que crescem quando o ambiente em questão são as redes sociais. Casos como “miiiiitooo boooooooooom” ou “kkkkkkkk” são vistos constantemente nestes ambientes, sendo necessária uma normalização prévia destes termos, como apresentado por (HAN; BALDWIN, 2011) para textos da língua inglesa. Esta abordagem foge do objetivo proposto, visto que a língua portuguesa não possui tantos trabalhos relacionados quanto à inglesa. Uma implementação simples para normalização da maioria dos jargões identificados foi realizada, visando melhorar o desempenho dos classificadores.

Com a Figura 3.4 é possível visualizar cada etapa do tratamento que será realizado na base textual, resultando na estruturação do texto.

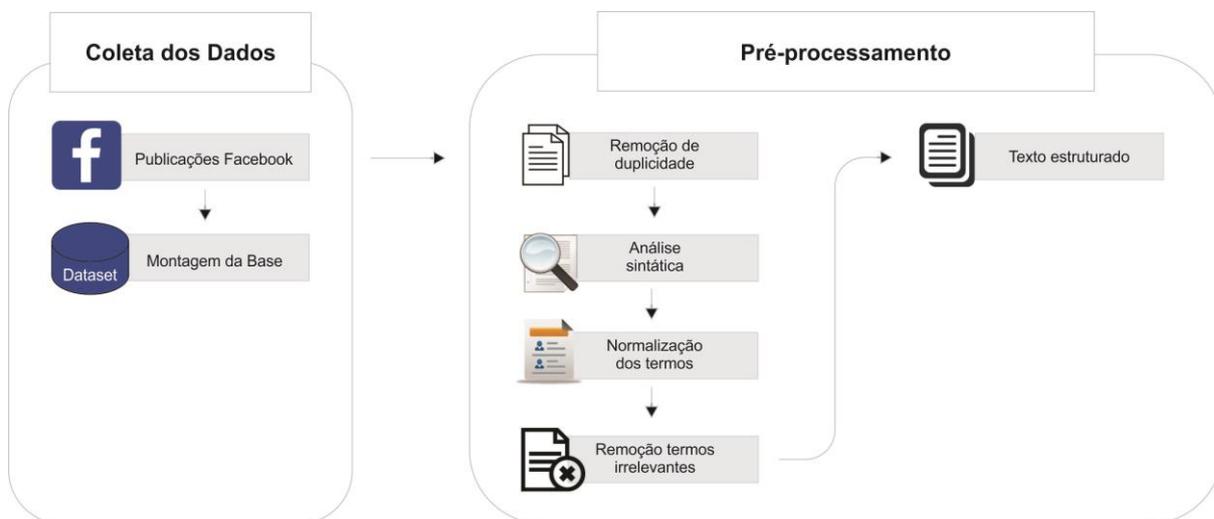


Figura 3.4: Etapas do pré-processamento.

Ainda na etapa do pré-processamento, foram eliminados todos os caracteres que não possuem composição com as palavras, isto é, caracteres que podem ser removidos sem perder a forma gramatical da palavra, como barras (/ ou \), operadores aritméticos (*, +, - e /), operadores relacionais (>, < e =), entre outros. Imagens e símbolos que representam uma expressão facial, caracterizada e conhecida como *emoticon*, também serão removidos do

texto. Há autores (GONÇALVES et al., 2013) que utilizam *emoticons* para caracterizar o estado emotivo da mensagem, entretanto, tais análises estão fora do escopo do presente trabalho.

A etapa referente à normalização morfológica, ao invés de eliminar os termos irrelevantes reduz os radicais dos termos, esse procedimento também é denominado como lematização ou *stemming* (FRAKES, 1992). Todavia, essa variação depende da tarefa de mineração de texto adotada. Na presente proposta, que utiliza léxicos, a redução de uma palavra ao seu radical afetará a comparação dos termos do texto com os termos dos léxicos.

Após as etapas citadas acima, os termos resultantes são submetidos ao léxico LIWC e aos léxicos especializados em termos afetivos, bem como à seleção dos termos com base em sua frequência, construindo assim, um vetor de características textuais para a inferência da personalidade. O LIWC será explorado na próxima subseção.

3.4.1. Utilização do Léxico LIWC

Conforme apresentação do segundo capítulo, a forma que as pessoas escrevem fornece abertura para observação de aspectos emocionais e cognitivos. Pesquisadores demonstram diversas evidências que sugerem que a saúde física e mental das pessoas está correlacionada com as palavras que elas usam (PENNEBAKER et al., 2007). As análises de textos, com base nesses estudos, indicam que os indivíduos tendem a usar determinadas classes de palavras, palavras de negação, palavras de afirmação, ou tipo de pontuação, entre outros.

Dessa maneira, para detectar a personalidade por meio de textos, são utilizadas informações relacionadas com a forma que um indivíduo escreve, observando a linguagem e as propriedades individuais das palavras presentes em seus textos. Para essa tarefa são utilizados recursos lexicais no presente trabalho.

As principais potencialidades de utilização de léxicos estão relacionadas com o número e a natureza dos atributos que descrevem, baseando-se em estatísticas sobre o uso de palavras individuais. Conforme observado na Seção 2.2, os estudos sobre a inferência de personalidade, utilizando léxicos e suas informações sobre os termos, podem ser utilizados como predador de personalidade, de maneira que cada uma das suas entradas, como adjetivos, discordância, palavras relacionadas a eventos psicológicos, dentre outros, tem a possibilidade de ser utilizada na construção de um modelo de inferência.

Dentre os recursos de léxicos utilizados nesse estudo, destaca-se o LIWC (*Linguistic Inquiry and Word Count*) que recentemente foi disponibilizado para português brasileiro (BALAGE FILHO et al., 2013), e possui cerca de 127.162 palavras catalogadas em 64 categorias, apresentadas anteriormente na Seção 2.2.1.

Com a utilização do acervo de palavras e as categorias contidas no LIWC, pretende-se obter um alto índice de características textuais acerca dos conteúdos. Este estudo é o pioneiro a utilizar o léxico LIWC na inferência de personalidade por meio de textos escritos em português.

Após o refinamento dos dados coletados nas redes sociais e o isolamento dos termos, é construída a base de treinamento do modelo de aprendizagem, contendo como atributos as categorias do léxico LIWC e como atributo-meta (classe) o resultado do questionário de personalidade preenchido pelo usuário. Primeiramente os dados são categorizados junto ao léxico LIWC, isto é, os termos tratados serão confrontados com o léxico e categorizados conforme as dimensões e categorias presentes no dicionário. Dessa maneira, se obterão características textuais que permitirão a análise de correlações com os traços de personalidade do indivíduo. Para ilustrar o processo, a Figura 3.5 apresenta a incorporação do LIWC na fase do pré-processamento.

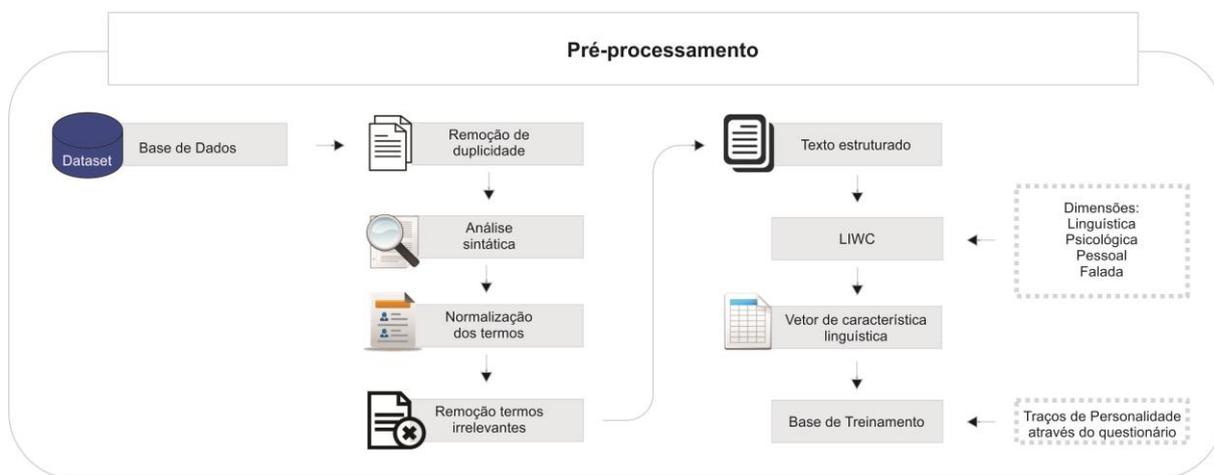


Figura 3.5: LIWC na etapa do pré-processamento.

Como sequência da evolução do método, a estrutura do vetor de característica é criada após os termos serem categorizados no léxico LIWC e realizadas as etapas iniciais do pré-processamento. Baseado no modelo de aplicação do léxico LIWC, feita por Mairesse (2005)

para a língua inglesa, o vetor contém o percentual correspondente ao número de palavras categorizadas no léxico e presente nas publicações dos usuários nas redes sociais. Ou seja, é computado o número de ocorrências em que a palavra é empregada nas publicações e posteriormente calculado o percentual em relação ao total de palavras publicadas.

Desta maneira, considerando a publicação “*Meu final de semana foi maravilhoso!!! Logo estarei de volta a esse estado maravilhoso.*”, exemplifica-se o processo de categorização das palavras no léxico, da seguinte forma: (i) a frase é submetida à etapa de pré-processamento, resultando uma lista de *tokens* “[*final*] [*semana*] [*maravilhoso*] [*logo*] [*estarei*] [*volta*] [*estado*] [*maravilhoso*]”, (ii) computa-se o total de palavras no texto (oito), (iii) a lista de *tokens* é submetida ao léxico, retornando apenas as palavras catalogadas e a quais categorias elas pertencem, (iv) a partir da lista de termos encontrados no léxico é sumariado o número de vezes que cada termo aparece, no exemplo, a palavra “*maravilhoso*” está presente duas vezes no texto, e (v) calcula-se o percentual de utilização de cada termo no texto, conforme Equação 3.1.

$$LIWC_{(t,d)} = \frac{(tf_{t,d} * 100)}{tt_d} \quad (3.1)$$

Na equação, $tf_{t,d}$ corresponde a ocorrência do termo t (“*maravilhoso*”) em um documento d . O valor tt_d é o total de palavras empregadas no documento d .

Desse modo, para a palavra “*maravilhoso*” obtém-se o valor 25. Esse valor é atribuído ao vetor de característica para a categoria do léxico em que o termo pertence. A Figura 3.6 ilustra o processo mencionado e o vetor de característica do LIWC no estudo.

As categorias do léxico LIWC utilizadas nesse experimento são apresentadas na Tabela 3.4. Ao conjunto padrão de 64 categorias, disponibilizado para o português brasileiro, foram acrescentadas cinco novas categorias, inspiradas no léxico LIWC de língua inglesa (PENNEBAKER et al., 2001). Essas categorias são independentes de idiomas, em razão de ser apenas contagem de características textuais, sendo tais categorias formadas pelo total da quantidade de: palavras, termos encontrados no léxico, palavras com mais de seis letras, palavras por frase e frases contidas no texto.

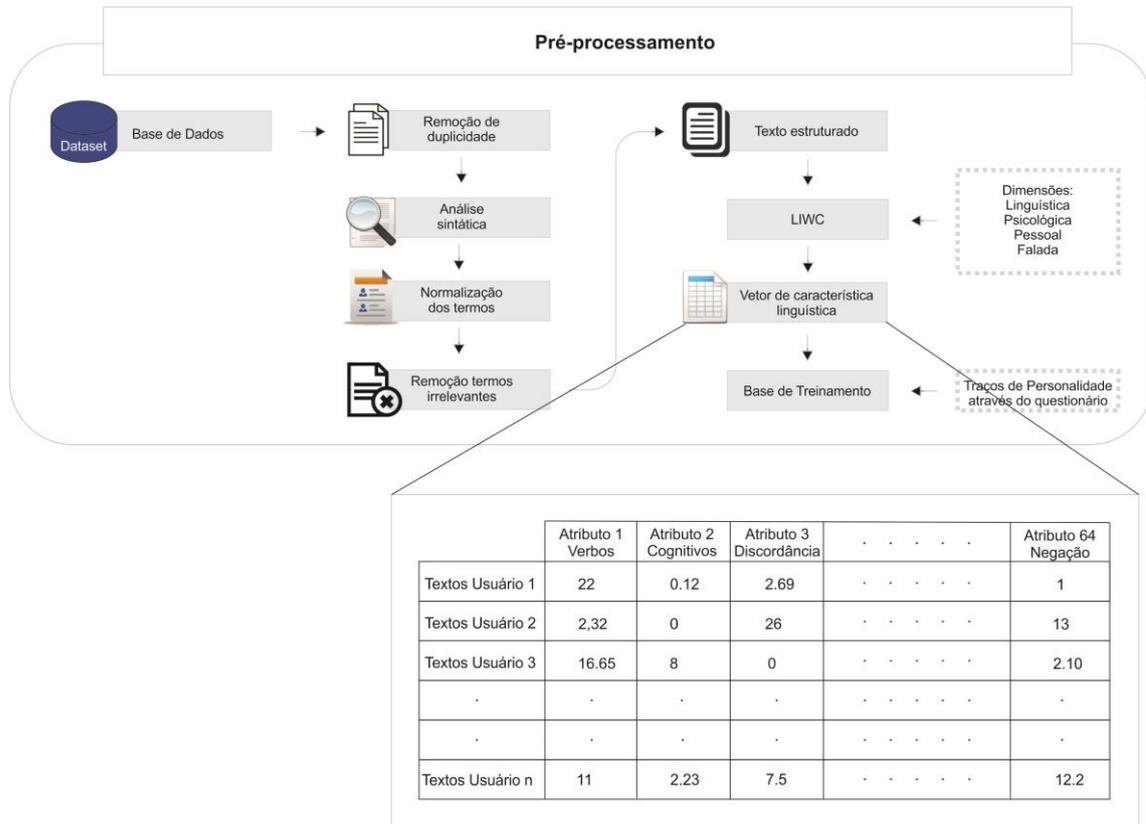


Figura 3.6: Demonstração do vetor de características.

Tabela 3.4: Categorias do LIWC para o português utilizadas no experimento.

Dimensões	Categorias
Linguística	total palavras (wc), palavras encontradas no léxico (dic), palavras com mais de seis letras (sixltr), palavras por frase (wps), total de frases (sentence) função de palavras (funct), total de pronomes (pronoun), pronome pessoal (ppron), 1ª pessoa do singular (i), 1ª pessoa do plural (we), 2ª pessoa (you), 3ª pessoa do singular (shehe), 3ª pessoa do plural (they), pronomes impessoais (ipron), artigos (article), verbos comuns (verb), verbos auxiliares (auxverb), passado (past), presente (present), futuro (future), advérbios (adverb), preposição (prep), conjunções (conj), negação (negate), quantificadores (quant), números (number), palavrões (swear)
Psicológica	social (social), família (family), amigos (friend), humano (human), afetivo (affect), emoções positivas (posemo), emoções negativas (negemo), ansiedade (anx), raiva (anger), tristeza (sad), cognitivos (cogmech), intuição (insight), causa (cause), discordância (discrep), tentativa (tentat), certeza (certain), inibição (inhib), inclusivo (incl), exclusivo (excl), perceptivo (percept), ver (see), ouvir (hear), sentir (feel), biológico (bio), corpo (body), saúde (health), sexual (sexual), ingestão (ingest), relatividade (relativ), movimento (motion), espaço (space), tempo (time)
Pessoal	Trabalho (work), conquista (achieve), lazer (leisure), dinheiro (money), religião (relig), morte (death)
Falada	Concordância (assent), sem fluência (nonflu), preenchimento (filler)

Para a tarefa de submissão das palavras contidas nas publicações no léxico LIWC, bem como a tarefa de pré-processamento do texto, foi criado um *software* capaz de computar cada categoria do léxico conforme descrito anteriormente. O *software* será descrito no próximo capítulo.

Observa-se que para construir um modelo de classificação com o intuito de categorizar os traços de personalidade de um indivíduo utilizando o *BigFive*, defronta-se com um problema de classificação multi-rótulo, visto que o *BigFive* possui cinco traços para uma pessoa. Entretanto, optou-se em decompor o problema multi-rótulo em cinco modelos pertencentes a cada traço, sendo eles: Extroversão, Neuroticismo, Socialização, Realização e Abertura à experiência.

Todavia, a base de dados do usuário e o vetor de característica linguística são utilizados para todos os modelos de treinamento, alterando apenas o atributo-meta conforme cada traço de personalidade. A Figura 3.7 ilustra a criação dos cinco classificadores e o vetor de característica com o atributo-meta rotulado.

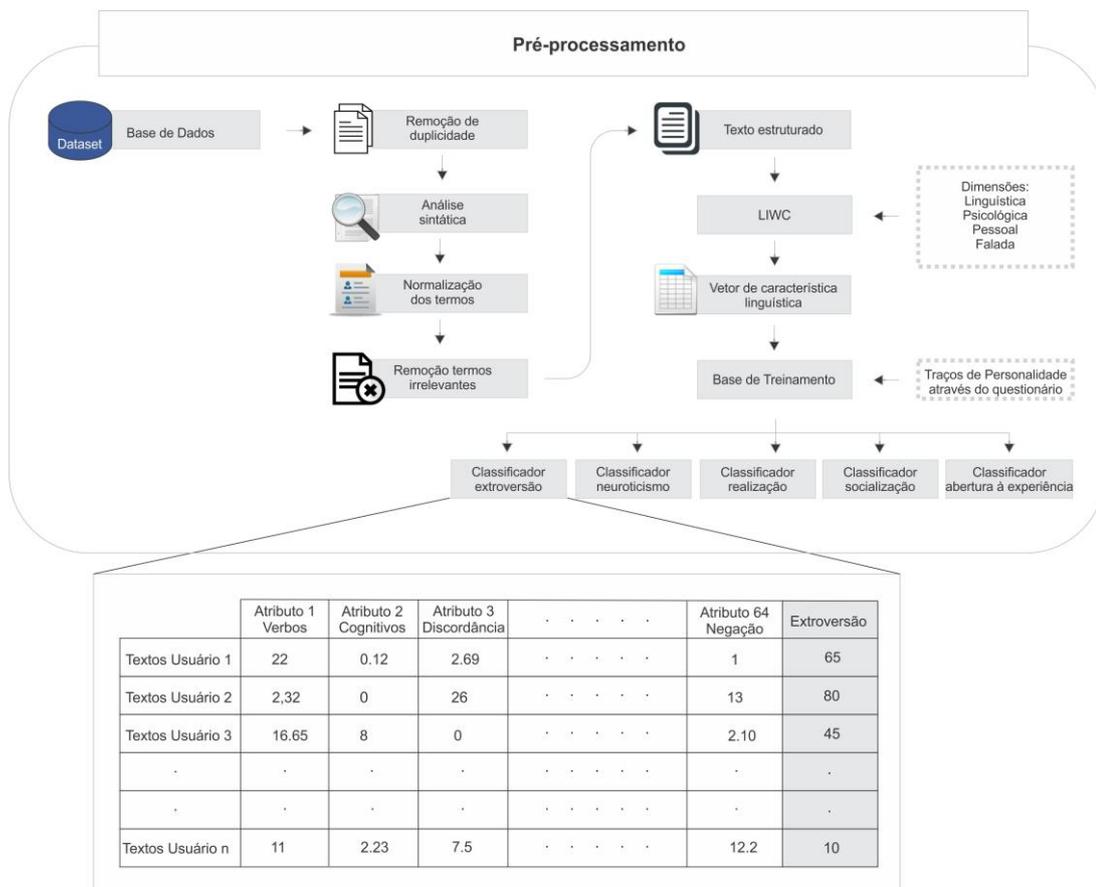


Figura 3.7: Estrutura da base de treinamento.

Até o momento a base de treinamento é criada a partir das características linguísticas extraídas dos textos com o auxílio do léxico LIWC. Essa estrutura de modelo permitirá comprovar a hipótese inicial desta pesquisa, que questiona se o reconhecimento da personalidade pode ser mensurado através de trechos escritos por seus autores. Todavia, presume-se que com o acréscimo de léxicos especializados em termos afetivos ao modelo, melhores resultados serão alcançados.

Na próxima seção serão demonstradas a extração da frequência dos termos no texto, e como esta será utilizada em conjunto com o modelo linguístico e afetivo.

3.4.2. Utilização de TF-IDF

Conforme apresentação do segundo capítulo há estudos de inferência de personalidade baseados na extração da frequência dos termos empregados em um documento. Dessa maneira, o modelo proposto para detectar a personalidade por meio de texto fará uso da representatividade que um termo tem para o seu conteúdo, através da abordagem TF-IDF (*Term Frequency - Inverse Document Frequency*), citada na Seção 2.2.

Com a utilização da abordagem TF-IDF, pretende-se obter uma lista de termos relevantes que em conjunto com o léxico LIWC e os léxicos afetivos, sejam capazes de se correlacionar com os traços de personalidade dos autores dos textos.

Para tal tarefa, os textos de cada usuário foram convertidos em uma matriz $m \times n$, na qual m corresponde ao documento de publicações do usuário e n corresponde a um termo, e conseqüentemente foi computado seu peso de frequência. Essa representação também é conhecida como *bag-of-words*, citada na Seção 2.2. A tabela 3.5 apresenta um exemplo da representação vetorial usando pesos TF-IDF no experimento.

Tabela 3.5: Representação das publicações usando peso TF-IDF.

Documentos	Termo 1	Termo 2	...	Termo n
Publicações do Usuário 1	0.15	0.30	...	0.20
Publicações do Usuário 2	0.20	0.10	...	0.90
...
Publicações do Usuário n	0.25	0.75	...	0.40

Para a extração dos termos com base em sua frequência (TF-IDF) e também para a criação da matriz, foi utilizada a ferramenta WEKA (WITTEN; FRANK, 2005), que conta

com todos os algoritmos e recursos adotados para essa tarefa. No próximo capítulo serão apresentadas mais informações sobre a ferramenta e a maneira com que os termos foram extraídos.

Os conjuntos de termos com base em sua frequência, extraídos por meio da abordagem TF-IDF, serão utilizados em conjunto com o léxico LIWC, com o objetivo de aumentar a eficácia do método para a inferência de personalidade a partir de texto. A Figura 3.8 exemplifica a união dos conjuntos de termos TF-IDF com o léxico LIWC na base de treinamento. Todavia, serão utilizados os mesmos conjuntos de termos para todos os modelos de treinamento, alterando apenas o atributo-meta conforme cada traço de personalidade.

Ainda, serão utilizados os conjuntos de termos frequentes, extraídos por meio da abordagem TF-IDF, como forma de seleção de palavras nos léxicos afetivos. Na próxima subseção, serão mostrados os léxicos afetivos utilizados em conjunto com o modelo linguístico e a maneira com que os conjuntos de termos frequentes serão aplicados a eles.

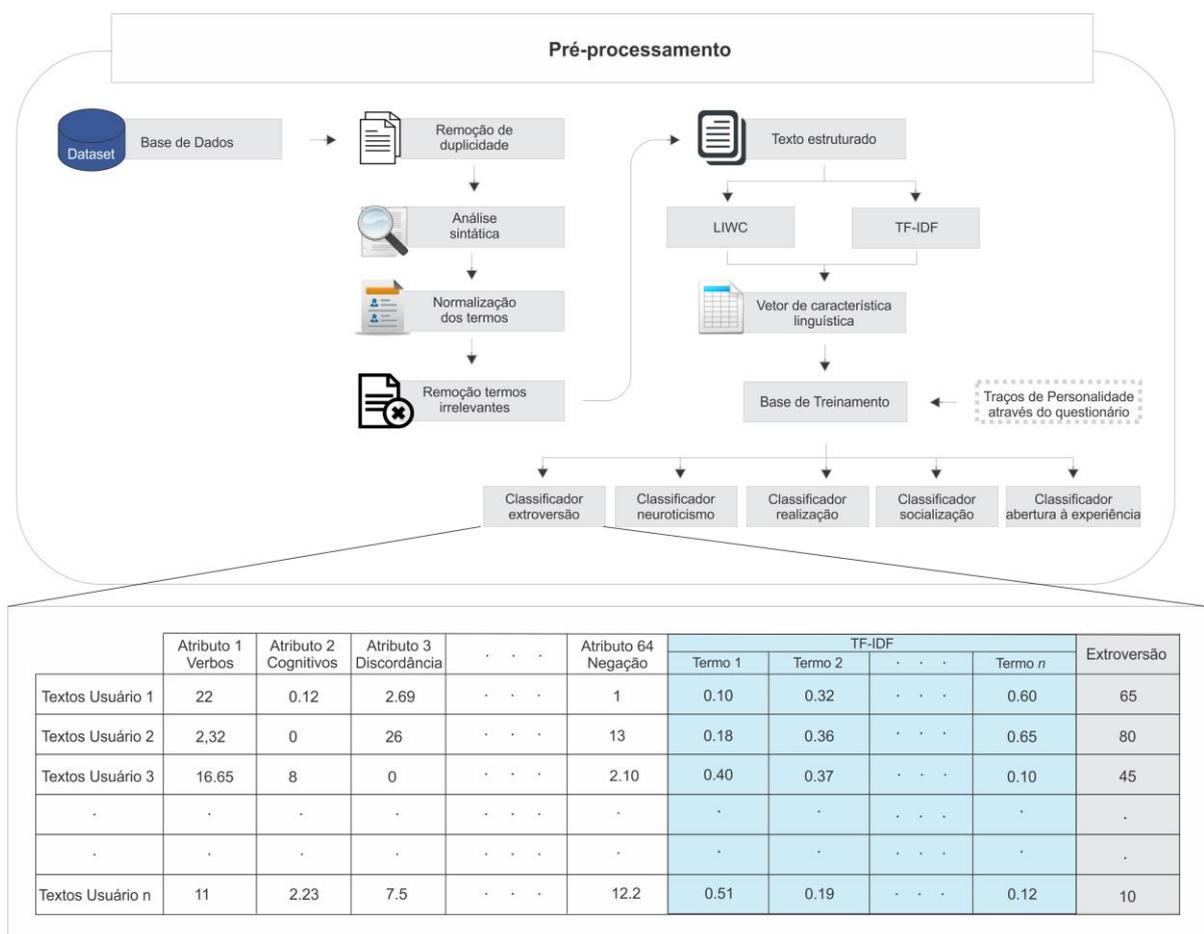


Figura 3.8: Estrutura da base de treinamento com os atributos LIWC e TF-IDF.

3.4.3. Utilização dos Léxicos Afetivos

A identificação de emoções contribui com a averiguação dos traços de personalidade, uma vez que a personalidade é identificada por características de um organismo autônomo que representa padrões consistentemente escolhidos de reação mental, incluindo emoções (MOFFAT, 1997). Vários pesquisadores mostraram que as emoções podem representar indicadores significativos para a descoberta da personalidade (Seção 2.2).

Da mesma forma, pretende-se utilizar características emocionais para o melhoramento da precisão do método proposto, por meio da utilização de léxicos afetivos construídos ou adaptados para a língua portuguesa por outros pesquisadores.

O estudo propõe, ainda, a união de vários léxicos afetivos, com o intuito de suprimir as deficiências de cada léxico no que diz respeito ao número de palavras catalogadas para a língua portuguesa e, também, identificar qual léxico afetivo possui melhor desempenho para o reconhecimento da personalidade.

Desse modo, serão analisadas duas abordagens para computar as palavras nos léxicos afetivos: (i) contabilizar as polaridades dos termos empregados nos textos (positivos, negativos e neutros) e, (ii) a partir da lista dos termos representativos nos textos, por meio da aplicação do método TF-IDF, submeter aos léxicos os termos considerando seu peso de relevância no texto e o peso dos léxicos para o termo, podendo esse último ser: positivo (1), negativo (-1) ou neutro (0). Esse cálculo é dado pela Equação 3.2, que atribui um valor $w_{t,d,l}$ para um termo t em um documento d para um léxico l . O valor $TFIDF(t,d)$ é o peso de frequência do termo t em um documento d , calculado pelo método TF-IDF, sendo que o valor P representa a polaridade do termo t no léxico l .

$$w_{t,d,l} = TFIDF(t,d) \times P_{t,l} \quad (3.2)$$

A seguir, serão citados os léxicos afetivos utilizados nesta pesquisa e como cada um deles é aplicado no método de inferência de personalidade:

SentiStrength. Possui uma base de dados de palavras em português catalogadas em uma escala de emoções positivas e negativas. Para cada entrada de texto, o *SentiStrength* classifica as palavras em um intervalo de 1 a 5 denotando palavras positivas, e -5 a -1 para palavras negativas. Por exemplo, a sentença “*Odeio o clima politico atual.*” aplicado ao *SentiStrength* é classificada da seguinte maneira: “*Odeio[-4] o clima politico atual.*”. A

palavra “*odeio*” é rotulada como negativa, obtendo a pontuação -4, que significa um forte sentimento negativo.

Para a primeira abordagem do experimento, que é computar o número de palavras positivas e negativas empregadas no texto, serão contabilizadas as palavras nas escalas positivas e negativas, ignorando seu peso de variação dos termos (intervalo de 1 a 5). Por exemplo, a sentença “*Odeio o clima politico atual.*” será estimada como uma emoção negativa presente na frase.

Deste modo, ao computar a quantidade de palavras positivas e negativas empregadas nos textos de cada indivíduo, será acrescentando ao vetor de características linguísticas a sumarização desses termos. A Figura 3.9 exemplifica a contabilização de termos positivos e negativos reportados por meio do *SentiStrength* e a adição dos atributos na base de treinamento.

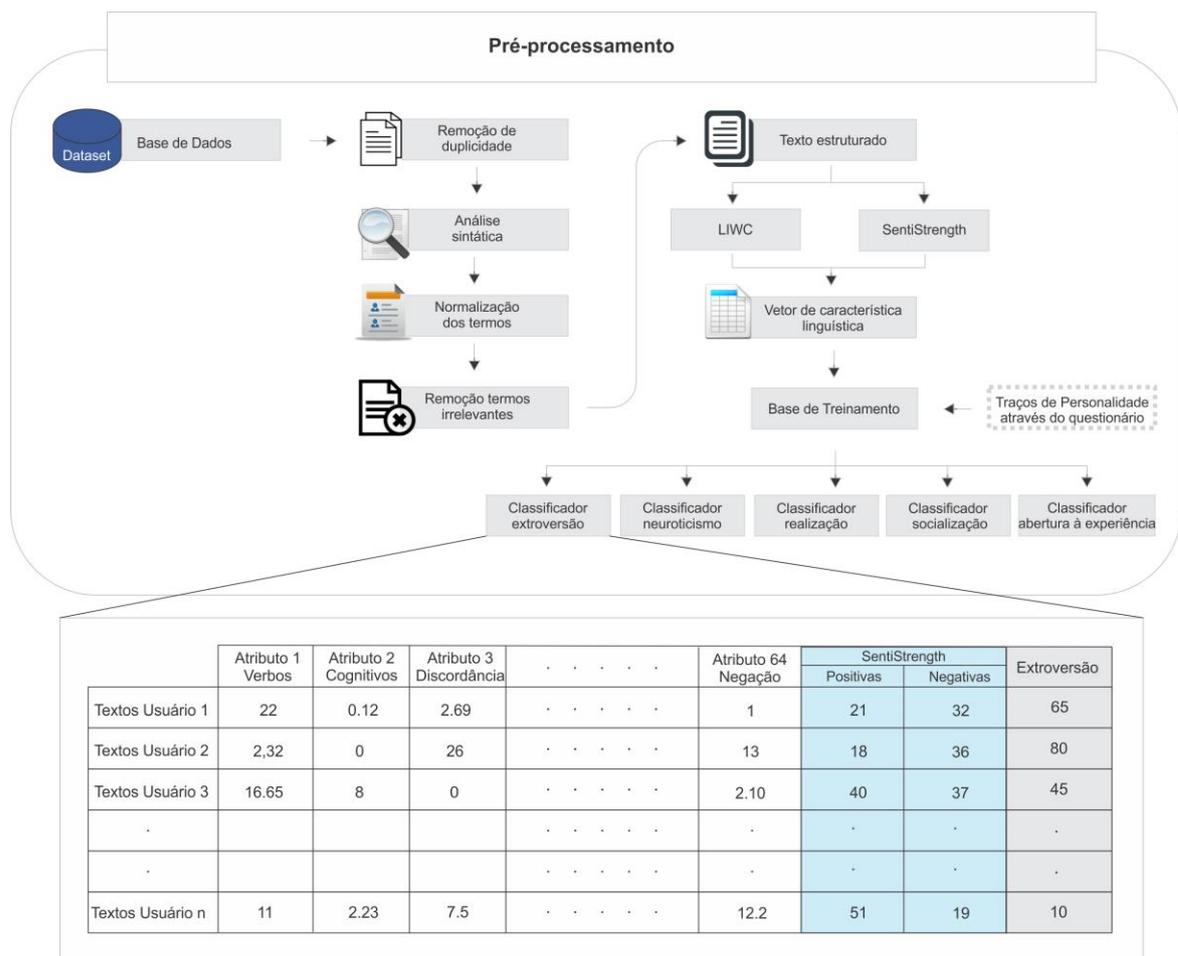


Figura 3.9: Estrutura da base de treinamento com a contabilização das emoções do léxico *SentiStrength*.

A partir dessa estrutura, será possível investigar o melhoramento na precisão do reconhecimento de personalidade, em decorrência da inclusão da contagem das emoções ao modelo.

Posteriormente, para a aplicação da segunda abordagem, um conjunto de termos será submetido ao léxico com seus pesos de frequências (TF-IDF), para averiguar quais desses termos estão contidos no léxico.

Para os termos presentes no léxico, será criado um vetor contendo o valor do peso do termo TF-IDF, multiplicado com peso resultante do léxico para tal termo, considerando as escalas de 1 a 5 para palavras positivas e -1 a -5 para palavras negativas. Por exemplo, dadas as seguintes palavras a partir do método TF-IDF: “*amor*” e “*tristeza*” com seus pesos de frequência 0.10 e 0.20 respectivamente. Ao confrontar a palavra “*amor*” no léxico *SentiStrength* é obtida a escala positiva de valor 3, que multiplicado com o peso 0.10, correspondente ao TF-IDF desse termo, resultará no valor 0.30 acrescentado ao vetor. Por sua vez, ao confrontar a palavra “*tristeza*” no léxico *SentiStrength* é obtida a escala negativa de valor -3, que multiplicado com o peso TF-IDF desse termo, resultará no valor -0,60 acrescentado ao vetor. O cálculo é dado pela Equação 3.2.

Deste modo, o vetor de termos do *SentiStrength* baseado na frequência do termo e no peso emotivo contido no léxico, será acrescentado ao vetor de características linguísticas, da mesma maneira que a primeira abordagem. A Figura 3.10 ilustra esse processo.

A partir da aplicação dessas abordagens será possível investigar o melhoramento na precisão do reconhecimento de personalidade em decorrência da inclusão do léxico *SentiStrength*, e ainda, diferenciar qual abordagem possibilitou tais melhorias.

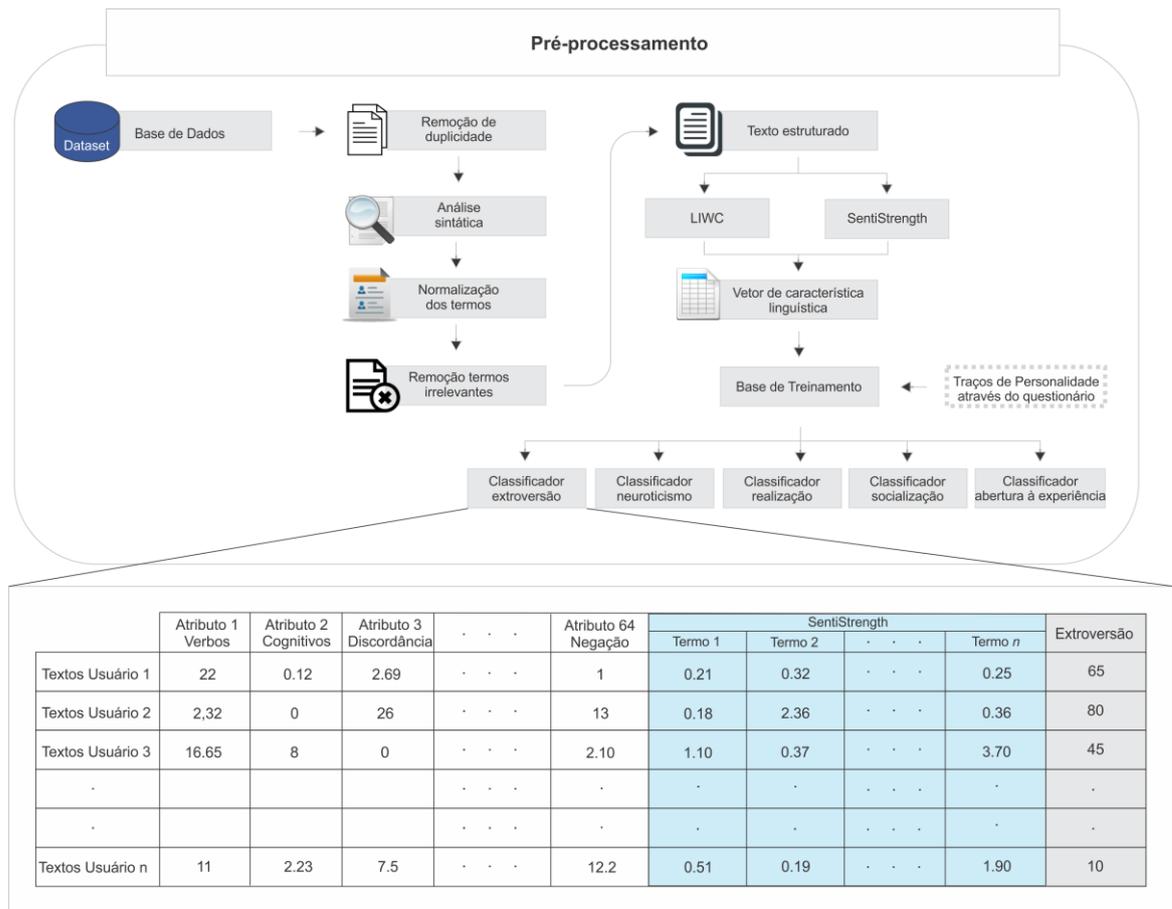


Figura 3.10: Estrutura da base de treinamento com a contabilização das emoções do léxico *SentiStrength* baseado na frequência dos termos.

AnewBr. O léxico possui uma lista de palavras afetivas que será utilizada para o confronto com o conjunto de termos da base textual criada, a fim de obter o nível emocional empregado nos textos.

A base do léxico possui 1.046 palavras para a língua portuguesa com valores de valência e alerta. Segundo Kristensen (KRISTENSEN *et al.* 2011) as emoções podem ser compostas por ao menos duas dimensões ortogonais, uma de valência (do desagradável ao agradável) e outra de alerta (do relaxado ao estimulado), em uma concepção definida como a teoria dimensional da emoção.

Dessa maneira, a primeira abordagem utilizando o léxico *AnewBr* será computar a média ponderada de valência e alerta empregada nos textos. A seguir são expostos detalhes sobre essa abordagem.

Os valores de valência e alerta estão compreendidos em um intervalo de 1 a 9. Sendo assim, palavras que possuem valores baixos, próximos de 1, por exemplo, apresentam

valência e alerta baixos, ou seja, desagradável e relaxado respectivamente. Já palavras com valores próximos de 9, apresentam valência e alerta altos, ou seja, agradável e estimulado respectivamente.

Por exemplo, a sentença “*Quero um mundo menos violento, cheio de paz!*” confrontada com a lista de palavras afetivas do *Anewbr*, é classificada da seguinte maneira: “*Quero um mundo*[Valência: 6,17 e Alerta: 5,01] *menos violento*[Valência: 1,46 e Alerta: 6,49], *cheio de paz*[Valência: 8,64 e Alerta: 3,71]!”. A partir dessa classificação é possível estimar a média ponderada de valência e alerta, conforme mostram as Equações 3.3 e 3.4. Elucida-se que M_v corresponde à média ponderada para a valência; M_A a média ponderada para o alerta; q_i a quantidade de vezes que uma palavra i é encontrada; v_i é o valor de valência de uma palavra i ; e A_i o valor de alerta de uma palavra i .

$$M_v = \frac{\sum_{i=1}^n (q_i * v_i)}{\sum_{i=1}^n q_i} \quad (3.3)$$

$$M_A = \frac{\sum_{i=1}^n (q_i * A_i)}{\sum_{i=1}^n q_i} \quad (3.4)$$

Deste modo, a média ponderada de valência e alerta empregada na frase “*Quero um mundo menos violento, cheio de paz!*” será de 5,42 e 4,89 respectivamente. Esses valores representam valência e alerta moderados.

Assim, pretende-se computar as médias ponderadas de valência e alerta empregada nos textos, acrescentando os valores ao vetor de características do LIWC.

Posteriormente, para a aplicação da segunda abordagem, será submetido ao léxico *AnewBr* um conjunto de termos com seus pesos de frequências (TF-IDF), para averiguar quais desses termos estão contidos no léxico.

Para os termos presentes no léxico, será criado um vetor contendo o valor do peso do termo TF-IDF, multiplicado com o peso resultante de valência e alerta do léxico para tal termo. Por exemplo, conforme apresentado anteriormente, a sentença “*Quero um mundo menos violento, cheio de paz!*” sendo confrontada com a lista de palavras afetivas do *Anewbr*,

é classificada da seguinte maneira: “*Quero um mundo[Valência: 6,17 e Alerta: 5,01] menos violento[Valência: 1,46 e Alerta: 6,49], cheio de paz[Valência: 8,64 e Alerta: 3,71]!*”. Para cada termo contabilizado no léxico, “*mundo*”, “*violento*” e “*paz*” serão criadas duas colunas no vetor, uma para armazenar o peso de valência e outra para o peso de alerta, contendo a multiplicação da frequência dos termos “*mundo*”, “*violento*” e “*paz*”. Para ilustrar o processo, a Figura 3.11 apresenta a estrutura da base de treinamento contendo o vetor do léxico *AnewBr* associado com o vetor de características do LIWC.

A partir da aplicação dessas abordagens, será possível investigar o melhoramento na precisão do reconhecimento de personalidade em decorrência da inclusão do léxico *AnewBr*, e ainda, comparar o *AnewBr* com os demais léxicos afetivos utilizados nos experimentos.

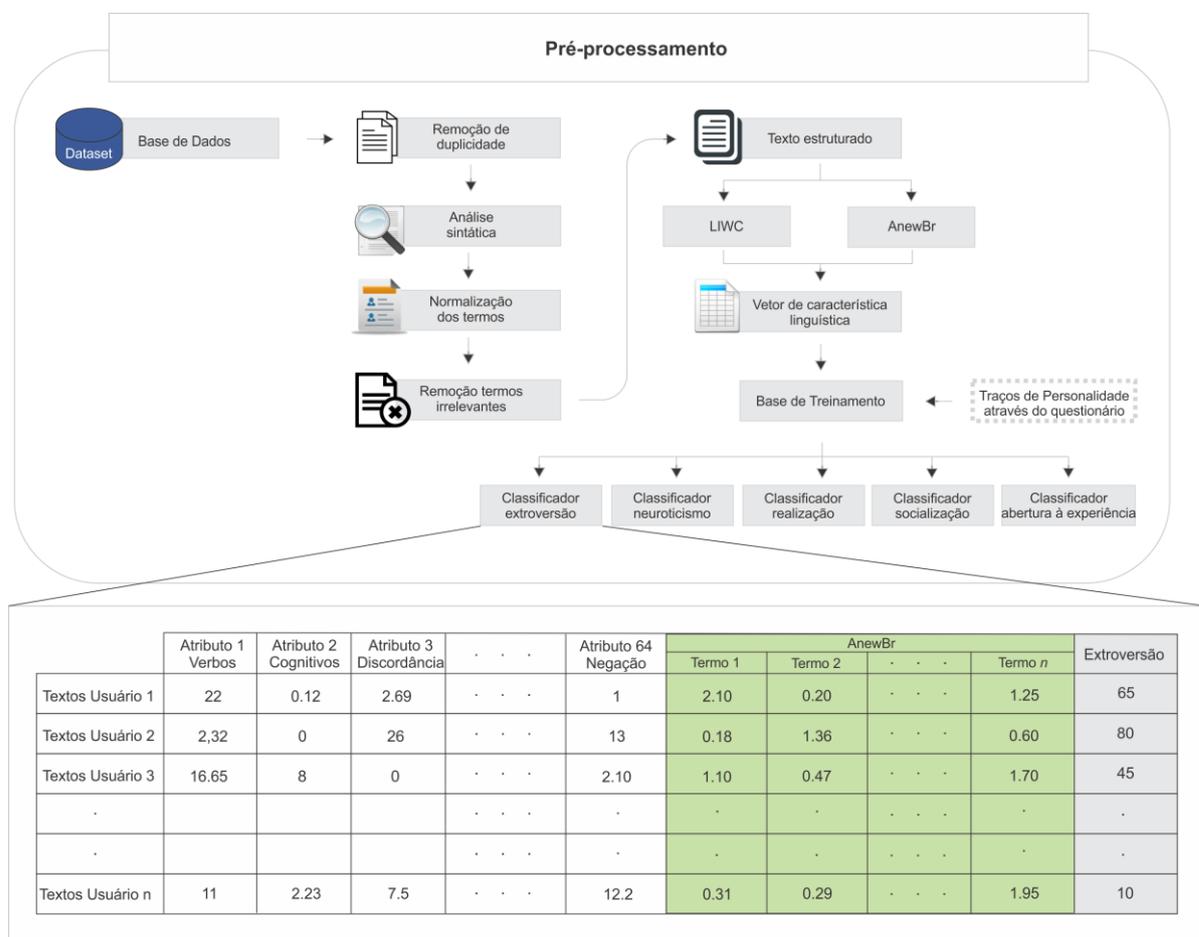


Figura 3.11: Estrutura da base de treinamento com a contabilização das emoções do léxico *AnewBr* baseado na frequência do termo.

SentiLex-PT. Dispõem de uma base de dados composta por 7.014 palavras em português catalogadas em polaridade emocional, a qual pode ser positiva, negativa ou neutra. O léxico foi concebido a partir de textos regidos em português de Portugal. Todavia, espera-se que o uso do léxico para textos escritos em português brasileiro possa obter grande cobertura de termos.

De maneira semelhante às abordagens aplicadas no léxico *SentiStrength*, pretende-se computar a quantidade total de palavras positivas, negativas e neutras resultantes do léxico *Sentilex-PT*. Pretende-se, ainda, aplicar um conjunto de termos TF-IDF ao léxico com o intuito de averiguar quais termos estão presentes, além de criar um vetor contendo o cálculo do peso da frequência com o peso do léxico para os termos, conforme Equação 3.2.

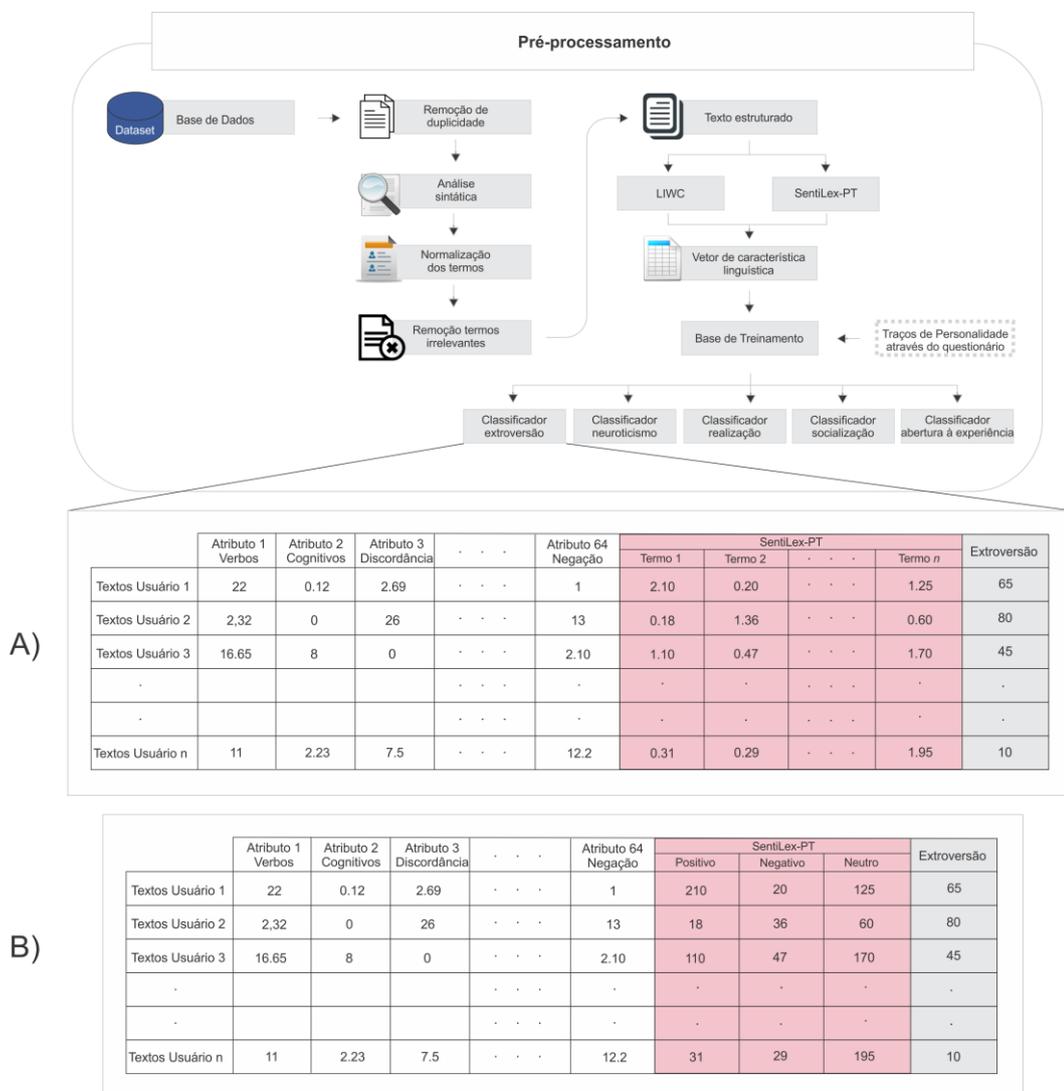


Figura 3.12: Estrutura da base de treinamento com a contabilização do *SentiLex-PT* em duas abordagens: (A) baseada na frequência dos termos; (B) soma dos termos emotivos.

Os valores de informações da polaridade associados aos termos no léxico são: -1 (um negativo) para palavras negativas, 0 (zero) para palavras neutras e 1 (um) para palavras positivas. Com o propósito de utilizar esses valores em conjunto com a frequência do termo, os valores de representação das polaridades foram alterados para: 1 (um) para palavras negativas, 2 (dois) para palavras neutras e 3 (três) para palavras positivas, mantendo a mesma escala de pontuação do termo. A Figura 3.12 apresenta a aplicação das duas abordagens no léxico *Sentilex-PT*.

OpLexicon. Léxico de sentimento criado para o idioma português brasileiro, constituído por cerca de 32.000 palavras (versão 3.0) anotadas com polaridades positivas (1), negativas (-1) e neutras (0).

Semelhante aos demais léxicos utilizados no experimento, serão contabilizados todos os termos positivos, negativos e neutros presentes nos textos de cada indivíduo, de maneira a agregar atributos no vetor de característica.

Além disso, os termos representativos no texto são submetidos ao léxico, criado através do método TF-IDF, com o objetivo de criar um vetor contendo o cálculo do peso da frequência com o peso do léxico para cada termo (Equação 3.2). Será alterada sua escala das polaridades, conforme modificação feita para o léxico *Sentilex-PT*. A Figura 3.13 apresenta a aplicação das duas abordagens no léxico *OpLexicon*.

Os léxicos afetivos podem ser utilizados de maneira conjunta, mesclando-se o seu uso em diversas combinações, com o intuito de maximizar a identificação das emoções empregadas nos textos. Aqueles que atribuem polaridade ao texto, também podem ser combinados entre si, pois cada léxico possui, na maioria das vezes, vocábulos diferentes. Por exemplo, o léxico *SentiLex-PT* possui polaridade de emoção para expressões idiomáticas da língua portuguesa, já o *OpLexicon* possui um número elevado de adjetivos e um conjunto de termos isolados, que o torna o léxico com maior número de termos catalogados.

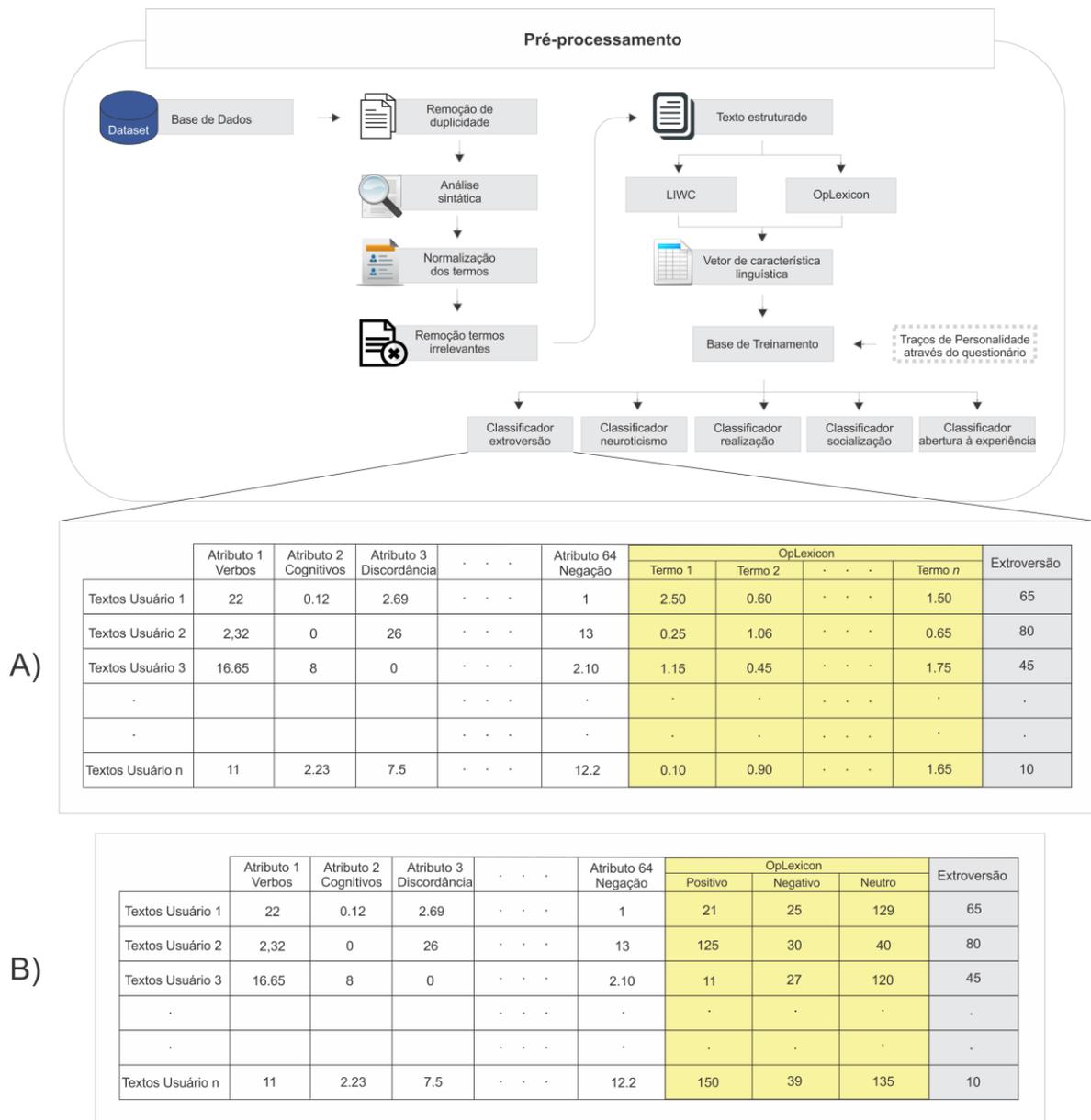


Figura 3.13: Estrutura da base de treinamento com a contabilização do *OpLexicon* em duas abordagens: (A) baseada na frequência dos termo; (B) soma dos termos emotivos.

A Figura 3.14 apresenta o vetor de características contendo a combinação dos termos positivos e negativos calculados do *OpLexicon* e as médias de estímulo de valência (desagradável ou agradável) e de alerta (relaxado ou estimulado) do *AnewBr*.

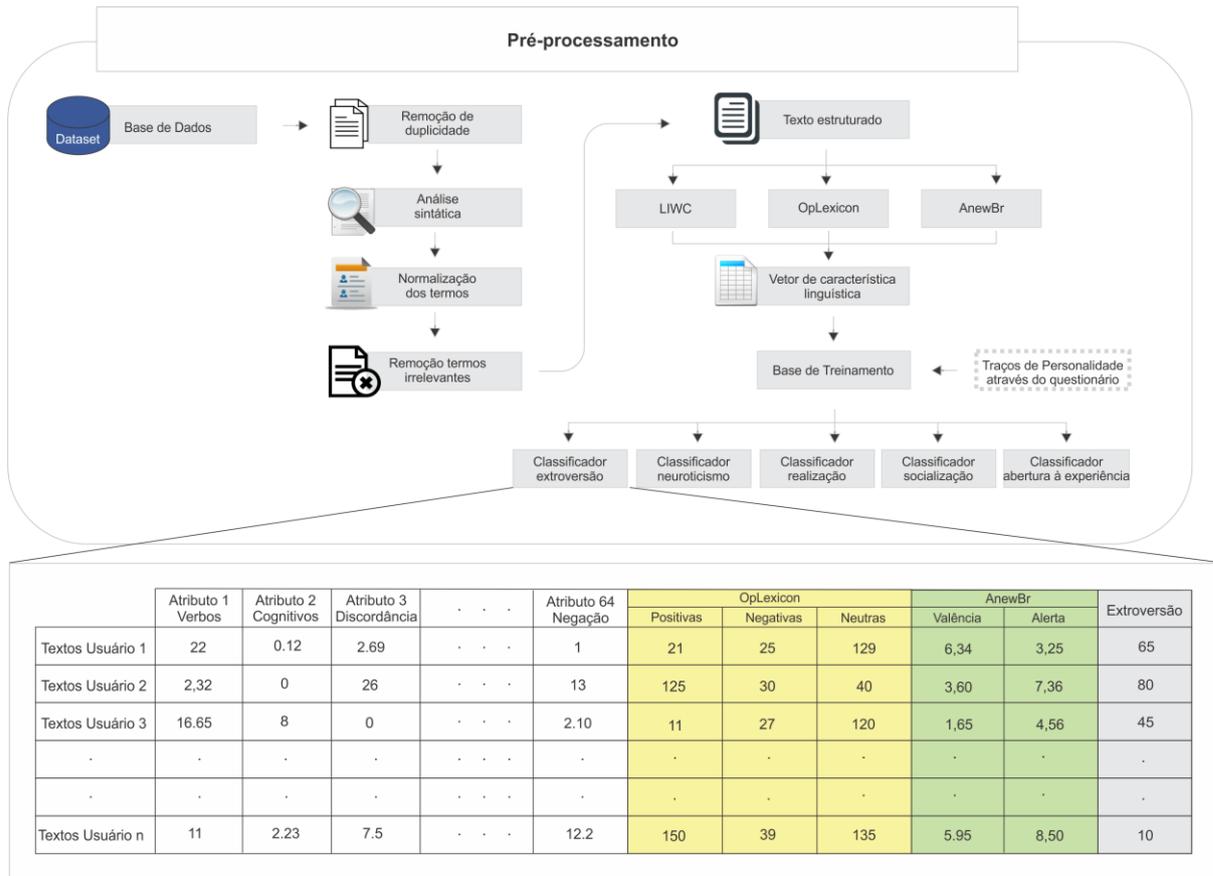


Figura 3.14: Demonstração da base de treinamento com a combinação dos léxicos: LIWC, *OpLexicon* e *AnewBr*.

Observando melhorias na precisão no reconhecimento da personalidade, em decorrência da inclusão dos léxicos afetivos, poderá ser analisado individualmente qual léxico possui melhor resultado. Ainda, tal processo permitirá comprovar a hipótese desta pesquisa, que questiona se com o acréscimo de um léxico afetivo é possível tornar o método de reconhecimento da personalidade robusto e com melhores resultados.

3.5. Considerações Finais

Nesse capítulo foram apresentadas as etapas do desenvolvimento de um método proposto para o reconhecimento da personalidade humana a partir de textos publicados *online*, escritos em língua portuguesa, por meio da análise de texto.

O modelo proposto possui como característica principal a utilização do léxico LIWC e léxicos especializados em termos afetivos para o reconhecimento da personalidade, tornando-se um diferencial em relação aos trabalhos na área para o português brasileiro.

Todavia, o método apresentado possui algumas limitações: (i) necessita de um *corpus* com textos rotulados para treinar os indutores; (ii) divide indutores em único-rótulo para resolver o problema multi-rótulo.

O próximo capítulo apresentará os detalhes de implementação, a descrição completa das ferramentas construídas para a coleta e pré-processamento dos dados e quais foram as ferramentas utilizadas para a execução do método.

Capítulo 4

Procedimentos Metodológicos

Neste capítulo são detalhados os procedimentos aplicados na construção e avaliação do método de inferência de personalidade em textos. O capítulo está dividido em três seções principais: a Seção 4.1, descreve as ferramentas utilizadas para a implementação do método descrito anteriormente; a Seção 4.2 descreve os algoritmos de aprendizagem de máquina utilizados para a construção de um modelo de indução; e a Seção 4.3 apresenta as métricas utilizadas para avaliar os algoritmos e o método.

4.1. Ferramentas de Software Utilizadas

Com a finalidade de implementar o método apresentado no capítulo anterior, foram desenvolvidas ferramentas que executam: (i) a extração dos traços de personalidade dos participantes, por meio de um inventário *online* (ii) a coleta dos textos publicados na rede social *Facebook*, (iii) o pré-processamento dos textos com o intuito de formatar e organizar de maneira adequada para a mineração dos dados, e (iv) a integração com a ferramenta WEKA para a mineração dos dados. A seguir são descritas as referidas ferramentas.

4.1.1. Aplicação do Inventário

Com a finalidade de extrair os traços de personalidade dos participantes, de maneira explícita, por meio de um inventário, foram apresentadas aos participantes as versões impressa (elaborada no processador de texto Microsoft Word) e *online* do questionário NEO-IPIP 210 para seu preenchimento.

Para a versão *online*, foi criada uma ferramenta eletrônica contendo todas as questões do questionário e a lógica de processamento dos resultados de cada traço de personalidade.

Ao final do preenchimento, o participante obteve sua pontuação e os resultados são armazenados em um banco de dados. A ferramenta apresentada foi implementada nas linguagens PHP e JavaScript (*JavaScript Object Notation*) utilizando a plataforma de desenvolvimento Eclipse.

Após a criação da ferramenta de inferência da personalidade por meio de inventário, também foi desenvolvido uma ferramenta para a tarefa de coleta de textos publicados na rede social dos participantes. O *software* será detalhado na próxima subseção.

4.1.2. Coleta de Textos

Com a finalidade de extrair as publicações dos usuários, foi desenvolvido um *software* capaz de obter as mensagens publicadas na rede social *Facebook*. Para esse trabalho, focou-se apenas o estudo em conteúdo textual, pois é nesse tipo de publicação que geralmente as pessoas expressam suas opiniões, sentimentos e exposição do seu comportamento acerca de um determinado assunto.

Na rede social *Facebook*, esse conteúdo está relacionado com o campo chamado “*status*”. Nesse campo a plataforma da rede social indaga o usuário com a seguinte mensagem “No que você está pensando?”, tal indagação, estimula os usuários à geração de conteúdo na rede social, podendo este ser: texto, imagem ou vídeo. O estudo desconsiderou o conteúdo gerado pelo usuário em formato de imagem ou vídeo.

Recentemente, a rede social *Facebook* efetuou mudanças em sua política de privacidade. Tal alteração resguarda o usuário à exposição de seu conteúdo fora da rede social, isto é, há restrição na extração de informações do campo “*status*” por meio de *software*, mesmo este ligado a sua *Application Programming Interface* (API).

A restrição mencionada pode ser revogada por meio de um processo minucioso de autorização e homologação do *software* por parte da equipe do *Facebook*, que autoriza o autor do *software* a extrair o conteúdo. Essa autorização está alienada com a conta do autor do *software* no *Facebook*, podendo esta ser invalidada a qualquer tempo.

Diante disso, o *software* deste estudo passou pelo processo rigoroso de autorização efetuado pela equipe do *Facebook*, e como resultado obteve autorização para extrair os textos publicados pelos usuários da rede. Contudo, para realizar essa extração, os usuários necessitaram autorizar, via aplicativo interno da rede, o acesso do *software* em seu perfil.

A autorização do *software* de extração de conteúdo apresentado ao usuário da rede social, também foi um fator na redução do volume de participantes do experimento, tendo em vista que nem todos os participantes visualizaram a notificação de autorização.

Todavia, os usuários que anuíram com a extração de suas publicações da rede social, tiveram todas as informações contidas em seu perfil extraídas pelo *software*, desde seu ingresso na rede social.

Para a extração das mensagens do *Facebook*, o *software* realizou a comunicação com a rede social através de troca de informações com a *Application Programming Interface* (API) do *Facebook*. As interações foram efetuadas por meio de solicitação HTTP (*Hypertext Transfer Protocol*) a partir do *software* que contém algoritmos escritos com capacidade de se comunicar com tal API. Uma vez que o *software* realiza essa comunicação, o mesmo estará apto para solicitar autorização e posteriormente coletar informações e publicações que um determinado usuário realizou na rede.

O *software* apresentado foi implementado na linguagem PHP, utilizando a plataforma de desenvolvimento Eclipse, seguindo os padrões de programação contidas no documento *Facebook for Developers*¹³ disponível na rede social.

Posteriormente à etapa de seleção e criação da base de dados, o processo avança para a etapa de pré-processamento, responsável pela limpeza e representação dos dados. Para essa tarefa foi desenvolvida uma ferramenta que será detalhada na próxima subseção.

4.1.3. Pré-Processamento

Para a etapa do pré-processamento dos dados, foi criado um *software* em linguagem *Delphi*, capaz de efetuar as etapas desse processo. A partir do fornecimento dos textos, o *software* realiza diversas tarefas de pré-processamento, sendo elas: elimina os sinais de pontuação, isola os termos (tokenização), converte as letras maiúsculas para minúsculas, normaliza jargões e vocábulos utilizados pelos usuários nas redes sociais, como por exemplo, as expressões “vc”, “amg” e “vdd” que correspondem aos termos “você”, “amigo” e “verdade” e, por fim, eliminam palavras de pouca importância para a representatividade do texto (*stopwords*).

Desse modo, o *software* representa o documento de forma vetorial utilizando o método *bag-of-words*, com o intuito de submeter as palavras ao léxico LIWC e aos léxicos afetivos.

¹³ <https://developers.facebook.com/docs/>

Após o processamento dos léxicos é exportado um arquivo, no formato ARFF (*Attribute-Relation File Format*), contendo os atributos e o atributo-meta, devidamente rotulado, para a etapa de mineração dos dados. Para essa fase foi utilizado um *software*, em linguagem Java, que possui integração com a ferramenta WEKA. Tal ferramenta será descrita na próxima subseção.

A Figura 4.1 ilustra a tela principal do *software*, responsável pela etapa de pré-processamento descrita.

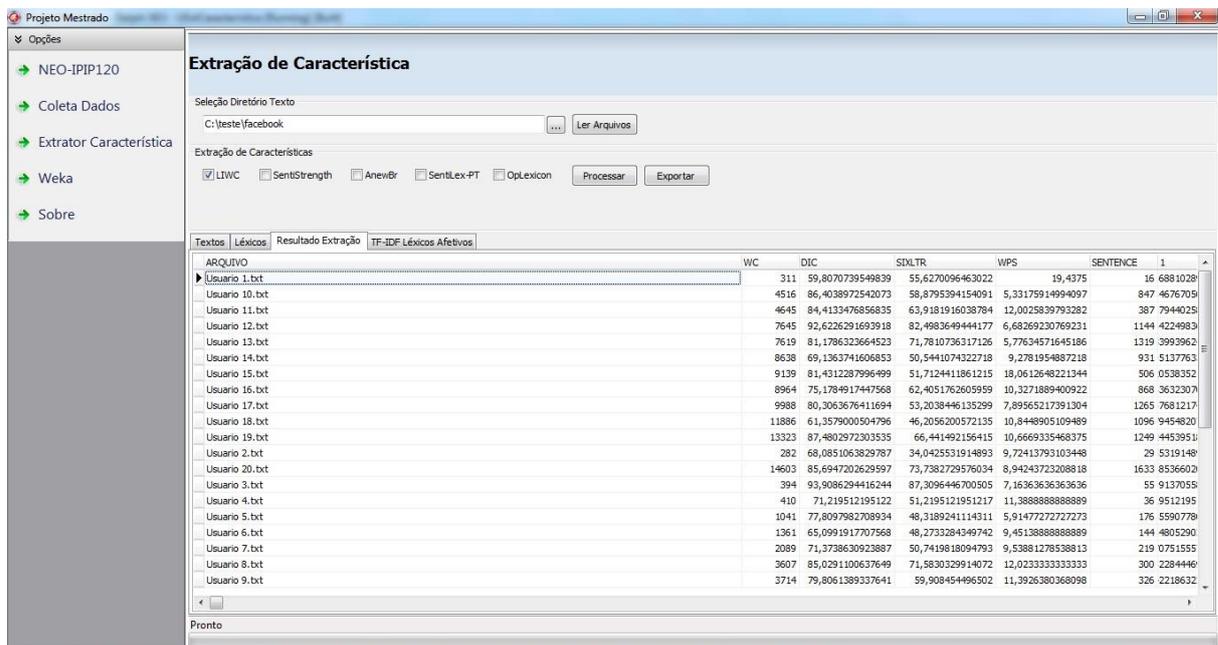


Figura 4.1: Tela principal do *software* responsável pela etapa de pré-processamento dos dados.

4.1.4. Waikato Environment for Knowledge Analysis - WEKA

Para a etapa de mineração dos dados foi utilizada a ferramenta WEKA¹⁴ (WITTEN; FRANK, 2005) de código aberto, que tem por finalidade agregar diversos algoritmos dedicados ao estudo de aprendizagem de máquina e mineração de dados. Contém diversos algoritmos para todas as etapas de mineração de dados, que são: pré processamento, classificação, regressão, agrupamento, associação e visualização. Todavia, para esse estudo, a ferramenta foi utilizada para o processamento dos algoritmos de regressão e a aplicação do método TF-IDF nos textos.

¹⁴ <http://www.cs.waikato.ac.nz/ml/weka/>

A plataforma WEKA fornece uma integração a qualquer aplicação Java. Tal integração permite a execução dos algoritmos e uma grande flexibilidade na manipulação dos dados, bem como maior facilidade de customização de todas as etapas presentes em um sistema de mineração de dados. Para a presente pesquisa, esse recurso de integração foi utilizado permitindo a execução dos experimentos de forma dinâmica.

Ainda, na plataforma foi utilizado o filtro *StringToWordVector*, considerando TF-IDF, como forma de ponderação dos termos, e para essa operação foi utilizada uma lista de *stopwords* em Português (apresentada no capítulo anterior). Tal processo resultou em uma base contendo os termos representativos no texto, com seu peso de frequência para cada documento. Foram criadas diversas bases (termo com sua frequência x documentos) variando o número de termo para cada matriz, a saber, os números de termos foram: 68, 150, 200, 250, 300, 500, 750, 1.000 e 1.500 termos. O objetivo de criar diversas bases é examinar qual quantidade de termos permitirá melhor desempenho na predição do traço de personalidade. Detalhes desse experimento serão tratados no próximo capítulo.

Os algoritmos de mineração de dados, utilizados na ferramenta WEKA, para criar o modelo de inferência de personalidade e validar o método, serão apresentados na próxima seção.

4.2. Algoritmos de Aprendizagem de Máquina Utilizados

Para a tarefa de construir um modelo de predição capaz de reconhecer os traços de personalidade de um indivíduo por meio de texto, um conjunto de algoritmos, candidatos para essa tarefa, devem ser analisados e testados. Dessa maneira, esta seção tem por objetivo apresentar os algoritmos de aprendizagem de máquina, que serão aplicados para um modelo de extração dos traços de personalidade dos dados textuais.

Todavia, o reconhecimento de personalidade pode ser visto como um problema de classificação multi-rótulo, em que cada rótulo consiste em um traço de personalidade. Ao todo, serão preditos cinco traços de personalidade para cada indivíduo. Desta forma, optou-se por desmembrar o modelo de classificação em cinco partes, utilizando os mesmos atributos e instâncias do vetor de características para cada modelo. Por conseguinte, esse desmembramento permitirá a predição de somente um traço de personalidade, isso dependerá da tarefa de mineração de texto adotada. Para a execução do método serão analisados todos os traços descritos pelo modelo *BigFive*.

A base de treinamento é composta pelas instâncias e atributos previamente rotulados (aprendizagem supervisionada), por meio do vetor de característica apresentado na Seção 3.4. Essa base de treinamento tem por objetivo construir um indutor que possa determinar corretamente o atributo meta de exemplos ainda não rotulados (base teste). Para a validação de cada algoritmo, foi usada a técnica de *k-fold cross-validation* (validação cruzada), em que o valor de *k* é igual a 10. Dessa forma, o conjunto de dados será dividido em dez partes iguais, utilizando nove partes para treinar e uma para testar. Mais informações sobre validação cruzada serão expostas na próxima seção.

Os valores possíveis para o atributo meta serão representados pelo nível (0 a 100) correspondente ao traço de personalidade do indivíduo, caracterizando, assim, um problema de regressão. Para a execução do método foram adotados diferentes modelos de indutores de regressão, a saber: regressão linear, com a utilização dos algoritmos *Linear Regression* (WITTEN; FRANK, 2005) e *SMOreg* (SHEVADE et al., 2000), árvore de decisão, com o uso do algoritmo *M5P* (WANG; WITTEN, 1997) e algoritmos *Lazy*, com a aplicação de *LWL* (ATKESON; MOORE; SCHAAL, 1997) e *IBK* (AHA; KIBLER; ALBERT, 1991).

Os experimentos foram realizados na ferramenta WEKA (descrita na seção anterior), que conta com todos os algoritmos adotados para este estudo. Os indutores possuem um conjunto de parâmetros que podem ser ajustados pelo usuário. A ferramenta traz valores *default* para os parâmetros de cada indutor. Para o treinamento, foram mantidos os valores *default*. Com exceção do parâmetro *KNN* (número de vizinhos) do *IBK*, que utilizou o valor *default* (1) e o valor 3, e o parâmetro *Kernel* do *SMOReg* que utilizou o valor *default* (*PolyKernel*) e *Puk*.

Para a estimativa e precisão dos algoritmos na etapa de mineração dos dados, serão utilizados métodos de avaliação apresentados na próxima seção.

4.3. Avaliação dos Resultados

Com o intuito de reconhecer a eficiência do modelo proposto, o método de avaliação dos resultados de reconhecimento de personalidade irá analisar a precisão (taxa de erros) dos algoritmos e correlação de Pearson (BUSSAB; MORETTIN, 1986) para os traços, determinando qual algoritmo e léxicos são mais indicados para o problema. A medida de avaliação será calculada sobre os exemplos pertencentes aos conjuntos de teste.

Para a estimativa de precisão dos algoritmos, serão utilizados os métodos de avaliação apresentados a seguir:

Validação cruzada (*k-fold*). Consiste em dividir o conjunto total de dados em k subconjuntos mutuamente exclusivos do mesmo tamanho. A partir disso, um subconjunto é utilizado para teste e os $k-1$ restantes são usados para estimação dos parâmetros e calcula-se a acurácia do modelo. Este processo é realizado k vezes alternando de forma circular o subconjunto de teste.

Correlação de Pearson. O coeficiente de correlação de Pearson é uma das medidas mais usadas para quantificar o grau de associação linear entre duas variáveis quantitativas. O coeficiente de correlação de Pearson (r), apresentado na Equação 4.1, quantifica a semelhança existente entre dois vetores de valores. Esse coeficiente sempre varia de -1 a 1. Uma relação diretamente proporcional exata ocorre quando o coeficiente é igual a 1. Uma relação inversa exata ocorre quando o coeficiente é igual a -1. Quando o coeficiente é igual a 0, significa que não existe relação linear entre os valores (WEINERT, 2010).

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (4.1)$$

Na equação, r representa o coeficiente de correlação, n o número de elementos no vetor, x o vetor de valores dos dados reais, y o vetor dos valores dos dados obtidos, i representa o i -ésimo elemento do vetor, \bar{x} a média dos valores do vetor x e \bar{y} a média dos vetores de y .

Para comprovarmos se o coeficiente de correlação é significativo, deverá ser realizado o seguinte teste de hipóteses: $H_0 : \rho = 0, H_1 : \rho \neq 0$. A estatística de teste é dada pela Equação 4.2, com $n-2$ graus de liberdade na tabela t de *Student* (n representa o tamanho da amostra). Caso o valor de t_c seja superior ao valor crítico de t , devemos rejeitar a hipótese nula. Se a hipótese nula, ao nível de significância α (por exemplo, 0,05) for rejeitada, podemos concluir que efetivamente existe uma relação significativa entre as variáveis.

$$t_c = r \cdot \sqrt{\frac{n-2}{1-r^2}} \sim t_{n-2;\alpha} \quad (4.2)$$

De acordo com (APPOLINÁRIO, 2006) o coeficiente de correlação de Pearson pode ser interpretado e avaliado qualitativamente conforme apresentado na Tabela 4.1.

Tabela 4.1: Valores de referência para a interpretação do coeficiente de correlação Pearson.

Adaptado de (APPOLINÁRIO, 2006).

Valores da Correlação	Interpretação
0,00	Nula
0,01 até 0,10	Muito Fraca
0,11 até 0,30	Fraca
0,31 até 0,59	Moderada
0,60 até 0,80	Forte
0,81 até 0,99	Muito Forte
1,00	Absoluta

Medida RMSE (Root Mean Squared Error). Consiste na raiz quadrada da média da diferença ao quadrado entre os valores reais e preditos para um atributo-meta (TORGO, 1995). Seja h_i a hipótese construída pelo algoritmo na i -ésima partição, n_{teste} o número de exemplos do arquivo de teste, y'_j corresponde ao valor predito pelo algoritmo no j -ésimo exemplo de teste e y_j é o valor real do atributo-meta deste mesmo exemplo. A equação é mostrada a seguir.

$$RMSE(h_i) = \sqrt{\frac{1}{n_{teste}} \sum_{j=1}^{n_{teste}} (y'_j - y_j)^2} \quad (4.3)$$

Teste Estatístico. Os testes estatísticos são uma estratégia para validação de resultados, que têm recebido crescente atenção da comunidade de aprendizagem de máquina. Seu principal objetivo é comparar os resultados obtidos em um conjunto de dados por diferentes algoritmos, com o intuito de verificar se esses diferem uns dos outros com relevância estatística.

Dessa maneira, para comparar a precisão entre os algoritmos, serão utilizados testes estatísticos para verificar se houve ou não melhorias no processo de previsão dos traços de personalidade. Para essa tarefa, será aplicado o teste de Friedman (FRIEDMAN, 1937) que trata-se de um teste estatístico não-paramétrico, cuja execução não requer o conhecimento da distribuição da variável da população. Esse teste permite ranquear os algoritmos utilizados nos experimentos e obter a resposta de qual possui o melhor comportamento. O principal objetivo do teste de Friedman é conferir se os classificadores gerados possuem diferenças expressivas, sendo caracterizada hipótese nula (quando não existem diferenças entre os algoritmos).

Para verificar se existe ou não correlação entre os dados, deve-se fazer o somatório das variâncias dos ranques, para então, através deste somatório calcular a probabilidade do valor ser superior ou igual à variância obtida utilizando a distribuição *qui-quadrada* com $k-1$ graus de liberdade. O resultado numérico final deste teste apresenta um nível de significância (p-valor). Quando tal valor for menor que 0.05, então é aconselhado rejeitar a hipótese nula. Em caso de rejeição da hipótese nula, pode-se prosseguir com um teste *post-hoc* de Nemenyi (CORDER; FOREMAN, 2011) para detectar quais são as diferenças entre os classificadores.

Esses testes permitirão afirmar, com alto fator de certeza, se os modelos gerados apresentam forte correlação com os traços quantificados pelo Inventário NEO IPIP-120. Também poderão comparar os modelos gerados com diferentes algoritmos de aprendizagem de máquina para a tarefa de regressão.

4.4. Considerações Finais

Nesse capítulo foram citadas as tecnologias utilizadas na implementação do método de reconhecimento de personalidade por meio de textos, regidos em língua portuguesa, bem como a ferramenta de mineração dos dados. Também foram apresentadas as métricas utilizadas para avaliar o desempenho do método.

O capítulo seguinte apresentará os principais experimentos que foram realizados com o método de inferência de personalidade.

Capítulo 5

Experimentos e Análise dos Resultados

Este capítulo descreve os experimentos realizados para a validação e avaliação de desempenho do método proposto no Capítulo 3. Neste contexto, a primeira seção apresenta a formação da base textual. Em seguida, são apresentados os resultados da aplicação do método TF-IDF e, posteriormente, os resultados dos léxicos. Dando continuidade à amostra dos resultados, é apresentada a combinação dos léxicos afetivos. Por fim, são avaliados os resultados obtidos para a tarefa de reconhecer a personalidade por meio de textos.

5.1. Formação da Base de Dados Textual

A base textual dessa pesquisa faz uso do volume de informações gerada na rede social *Facebook*. O processo tem como objetivo capturar as mensagens publicadas pelos 256 participantes (os mesmos participantes da seção anterior), permitindo medir suas dimensões de personalidade por meio dos conteúdos postados.

A base de dados é composta por textos publicados no *Facebook*, independentemente da data em que o participante postou a mensagem. Dessa forma uma publicação realizada há três anos pode ser significativa na construção do modelo.

A Tabela 5.1 apresenta uma amostra de textos que compõem a base de dados, coletadas dos participantes do experimento.

Tabela 5.1: Amostra dos textos coletados na rede social *Facebook*.

Usuário	Texto publicado
Usuário 1	Estou feliz por ter conquistado mais uma
Usuário 1	E claro parabéns para nossos verdadeiros guerreiros que não deixaram de acreditar em nenhum momento! Parabéns @atletico É CAMPEÃO
Usuário 2	Só quero o que me faz completamente feliz!
Usuário 2	Cidadãos brasileiros estão sendo atacados por delinquentes na rua #vergonha @PCERJ
Usuário 2	Eu talvez não tenha muitos amigos. Mas os que eu tenho são os melhores que alguém poderia ter
Usuário 3	O sucesso é ir, de fracasso em fracasso, sem perder o entusiasmo !. #ToNaLutaBrasil

Com base nos dados coletados, se obteve um conjunto textual composto por 2.590.034 palavras, não exclusivas, presentes em 187.488 publicações realizadas pelos usuários participantes do experimento na rede social. A Tabela 5.2 apresenta outras informações da base textual coletada.

Tabela 5.2: Informações sobre o conjunto de dados.

Informações	Total
Usuários	256
Publicações	187.488
Média de publicações por usuário	732
Palavras	2.590.034
Média de palavras por usuário	10.117

Após a etapa de seleção e criação da base de dados, o processo avançou para a etapa de pré-processamento, executando todas as fases exploradas na Seção 3.4. Como resultado do pré-processamento obteve-se a representação dos dados em formato vetorial, permitindo-se a exploração dos léxicos e, conseqüentemente, a criação das bases de treinamento a serem aplicadas aos algoritmos de aprendizagem de máquina. A seguir, a próxima seção explora os experimentos e resultados da inferência de personalidade, utilizando o método TF-IDF.

5.2. Reconhecimento de Traços com TF-IDF

Para validar o reconhecimento de personalidade por meio de texto, com a utilização da abordagem TF-IDF, obteve-se uma lista de termos relevantes que são capazes de se correlacionar com os traços de personalidade de seus autores.

Por consequência, os textos de cada usuário foram convertidos em uma matriz $m \times n$, na qual m corresponde à quantidade de usuários da base e n corresponde à quantidade de termos. Dessa forma, cada linha da tabela corresponde a um documento formado pelas publicações de um usuário específico e cada coluna corresponde ao peso do termo para os usuários. Foi utilizada a ferramenta WEKA (WITTEN; FRANK, 2005) para a extração dos termos com base em sua frequência (TF-IDF) e a criação da matriz, que conta com todos os algoritmos e recursos adotados para esta tarefa. Esse modelo de representação e a ferramenta foram discutidos na Seção 3.4.2 e no Capítulo 4 respectivamente.

Deste modo, com o propósito de determinar o melhor conjunto que se correlaciona com o traço de personalidade (atributo-meta), optou-se em extrair inúmeros conjuntos, por meio da ferramenta WEKA, variando o número de termos para cada conjunto. Nesse contexto, a Tabela 5.3 apresenta o total de conjuntos extraídos e a quantidade de termos para cada conjunto.

Tabela 5.3: Conjunto de termos TF-IDF extraídos dos textos.

Nome	Quantidade de termos
TF-IDF-68	68
TF-IDF-150	150
TF-IDF-200	200
TF-IDF-250	250
TF-IDF-300	300
TF-IDF-500	500
TF-IDF-750	750
TF-IDF-1000	1000
TF-IDF-1500	1500

A partir desses conjuntos foi possível obter uma lista de termos representativos das publicações dos usuários. Essa lista auxiliará na distinção da ocorrência de algumas palavras

que geralmente são mais comuns que outras nas publicações dos usuários. A seguir, serão listados os termos que fazem parte do conjunto TF-IDF-68, que contém 68 termos.

Tabela 5.4: Lista de 68 termos extraídos por meio de TF-IDF nas publicações dos usuários.

Termos
acho, amanhã, amigos, amo, amor, ano, apenas, aqui, assim, bem, boa, bom, brasil, cara, casa, coração, dar, deus, dia, dizer, então, estar, falta, faz, feliz, fica, ficar, fim, gente, hoje, hora, ir, lado, linda, mãe, melhor, menos, mim, mundo, nada, noite, novo, olha, onde, pai, parabéns, pessoa, pessoas, preciso, quer, quero, sabe, saudade, sei, semana, senhor, sim, tão, tarde, tempo, vai, vamos, vem, verdade, vida, você, vocês, vou

A lista de *stopwords* utilizada no processo de busca dos termos frequentes foi modelada para não eliminar palavras que estão contidas nos léxicos, pois se pretende utilizar tais listas para a seleção de palavras nos léxicos afetivos. Tal experimento será discutido adiante, na Seção 5.4.

Após a construção dos conjuntos de termos representativos no texto, os registros foram utilizados para formar as bases de treinamento e teste (de cada conjunto) à submissão aos algoritmos de regressão (citados na Seção 4.2). A avaliação do desempenho dos algoritmos foi realizada pelo método Validação Cruzada (*10-folds*). As próximas tabelas apresentam a correlação Pearson obtida entre cada traço de personalidade e os conjuntos TF-IDF criados.

Tabela 5.5: Resultados de correlação de Pearson do experimento TF-IDF para o traço Extroversão (* correlação significativa ao nível de 0,05).

Algoritmos	Extroversão								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
M5P	0.215*	0.134*	0.183*	0.133*	0.169*	0.143*	0.174*	0.168*	0.144*
LinearRegression	0.174*	-0.0008	-0.028	-0.050	0.042	0.066	0.184*	0.213*	0.264*
SMOReg	0.208*	-0.028	0.064	-0.017	0.042	0.076	0.182	0.196*	0.263*
SMOReg (Kernel=Puk)	0.095	0.087	0.089	0.074	0.077	0.043	0.043	0.024	0.001
LWL	0.017	0.078	0.018	-0.002	0.109	0.110	0.181*	0.168*	0.164*
IBK (KNN=1)	0.021	0.018	0.042	0.000	0.122	0.107	0.102	0.056	0.079
IBK (KNN=3)	0.016	0.043	0.013	0.001	0.011	0.053	0.070	0.145*	0.157*

Tabela 5.6: Resultados de correlação de Pearson do experimento TF-IDF para o traço Neuroticismo (* correlação significativa ao nível de 0,05).

Algoritmos	Neuroticismo								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
M5P	0.082	0.032	-0.018	0.013	0.012	0.142*	0.150*	0.044	0.011
LinearRegression	0.062	0.046	0.029	-0.026	0.095	0.107	0.012	0.009	-0.038
SMOReg	0.008	-0.022	0.056	0.062	0.091	0.107	0.013	0.007	-0.03
SMOReg (Kernel=Puk)	0.027	-0.038	-0.044	-0.074	-0.078	-0.027	-0.138	-0.143	-0.151
LWL	0.162*	0.137*	0.129*	0.114	0.074	0.103	0.086	0.036	0.047
IBK (KNN=1)	0.026	-0.041	-0.007	0.015	-0.005	0.039	0.001	0.024	-0.062
IBK (KNN=3)	0.097	0.089	0.067	0.052	0.083	0.058	0.020	0.020	-0.044

Tabela 5.7: Resultados de correlação de Pearson do experimento TF-IDF para o traço Realização (* correlação significativa ao nível de 0,05).

Algoritmos	Realização								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
M5P	0.054	-0.108	0.093	0.048	0.104	0.039	0.077	0.030	-0.008
LinearRegression	0.020	-0.071	0.037	-0.101	0.037	0.018	0.004	0.058	0.095
SMOReg	-0.020	0.030	-0.047	-0.021	0.022	0.003	-0.004	0.044	0.090
SMOReg (Kernel=Puk)	-0.007	-0.030	-0.036	-0.035	-0.019	-0.040	-0.049	-0.054	-0.066
LWL	0.079	0.048	-0.027	-0.0003	-0.060	0.020	0.011	0.006	0.040
IBK (KNN=1)	-0.046	-0.041	-0.015	0.017	0.103	0.112	0.078	0.051	0.047
IBK (KNN=3)	-0.014	-0.066	-0.081	-0.016	-0.023	0.056	0.079	0.041	0.095

Tabela 5.8: Resultados de correlação de Pearson do experimento TF-IDF para o traço Socialização (* correlação significativa ao nível de 0,05).

Algoritmos	Socialização								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
M5P	-0.002	-0.048	0.009	-0.055	-0.021	-0.041	0.015	0.062	0.064
LinearRegression	0.112	-0.022	0.013	0.067	0.067	-0.015	0.0006	0.026	0.036
SMOReg	-0.036	-0.074	0.082	0.049	0.053	-0.020	-0.013	0.029	0.023
SMOReg (Kernel=Puk)	0.034	0.034	0.047	0.053	0.056	0.057	0.045	0.038	0.026
LWL	-0.070	0.023	0.038	0.022	0.123*	0.058	0.091	0.109	0.085
IBK (KNN=1)	-0.039	-0.083	-0.016	0.026	0.026	-0.001	-0.050	0.032	0.077
IBK (KNN=3)	-0.103	-0.040	-0.048	-0.013	-0.011	0.011	0.092	-0.033	0.096

Tabela 5.9: Resultados de correlação de Pearson do experimento TF-IDF para o traço Abertura (* correlação significativa ao nível de 0,05).

Algoritmos	Abertura								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
M5P	0.030	0.059	0.042	0.102	0.08	0.070	0.072	0.096	0.133*
LinearRegression	0.111	-0.022	-0.000	-0.028	0.023	0.052	0.103	0.167*	0.181*
SMOReg	0.090	-0.034	-0.082	0.028	0.018	0.061	0.103	0.171*	0.172*
SMOReg(Kernel=Puk)	0.054	0.108	0.125*	0.129*	0.124*	0.117	0.112	0.103	0.087
LWL	0.034	-0.013	0.014	0.082	0.023	0.066	0.016	0.015	0.036
IBK (KNN=1)	-0.088	0.099	0.086	0.185*	0.143*	0.059	0.118	0.144*	0.172*
IBK (KNN=3)	-0.014	0.04	0.061	0.068	0.088	0.125	0.149*	0.143*	0.150*

Observa-se que para todos os traços de personalidade obteve-se uma correlação fraca, resultados entre 0,11 a 0,30, para os conjuntos de termos TF-IDF. Todavia, os resultados obtidos foram promissores, mesmo havendo uma baixa correlação, pois os conjuntos de frequência dos termos podem ser associados a outras abordagens de inferência de personalidade, visando melhores resultados. Essa associação será explorada em outros experimentos.

Cabe ainda detalhar os melhores resultados apresentados para cada traço, utilizando o método TF-IDF. Nota-se que os traços Extroversão, seguidos de Abertura à experiência e Neuroticismo, possuem uma correlação maior com as frequências dos termos comparados com os traços Socialização e Realização. A Tabela 5.10, apresenta os algoritmos e conjuntos de termos que atingiram maiores correlações com os traços de personalidade.

Tabela 5.10: Melhores correlações TF-IDF com os traços de personalidade.

Traço	Algoritmo	Valor de Correlação	Conjunto TF-IDF
Extroversão	Linear Regression	0,264	TF-IDF-1500
Neuroticismo	LWL	0,162	TF-IDF-68
Realização	IBK(KNN=1)	0,112	TF-IDF-500
Socialização	LWL	0,123	TF-IDF-300
Abertura	IBK(KNN=1)	0,185	TF-IDF-250

Nota-se, ainda, que o resultado obtido pelos algoritmos para os traços Neuroticismo, Realização, Socialização e Abertura, rejeitam a hipótese de que quanto maior o número de

termos, maior é a correlação com o atributo-meta. Evidencia-se que para tais traços o uso de determinados termos é frequente, que os usuários utilizam-se de vocábulos limitados e, ainda, que possuem a tendência de reutilizar palavras em seus discursos.

Para determinar a existência ou não de diferença estatística significativa entre os algoritmos, com todos os conjuntos do experimento, utilizou-se o teste de Friedman e *post-hoc* de Nemenyi para analisar quais algoritmos e/ou grupos de algoritmos são diferentes, ambos utilizando um nível de significância $\alpha = 0,05$.

Dessa maneira, para os traços Extroversão, Realização e Abertura, o teste de Friedman apontou que não existe diferença entre os algoritmos. Para os traços Neuroticismo e Socialização o teste apontou que existe diferença entre os algoritmos, cujos resultados do teste de Nemenyi podem ser observados na Figura 5.1.

Observa-se que para o traço Neuroticismo (Figura 5.1a) o algoritmo SMOReg (*Kernel=Puk*) possui o maior *rank* médio, o que caracteriza o pior caso. Dessa forma, os algoritmos LWL, M5P e IBK(*KNN=3*) são significativamente melhores que IBK (*KNN=3*) e SMOReg (*Kernel=Puk*). Ainda, para o traço Socialização, observa-se que os algoritmos estão todos conectados, o que significa que não há diferenças estatísticas significativas entre eles, destacando-se apenas o algoritmo LWL com menor *rank* médio.

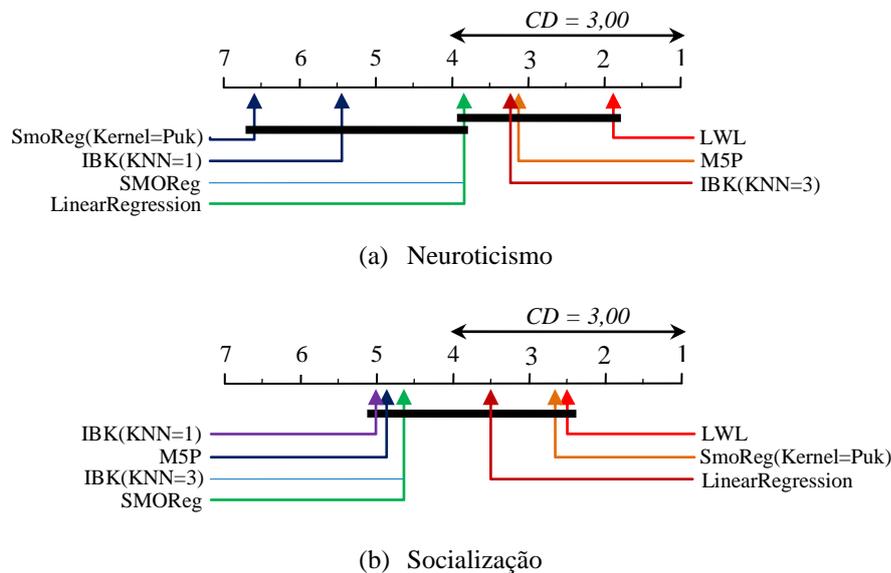


Figura 5.1: Resultado do teste de Nemenyi para o comparativo dos algoritmos do experimento TF-IDF para os traços: Neuroticismo e Socialização.

As bases de dados desse experimento foram normalizadas [0,1] com o intuito de obter melhores resultados do que os apresentados anteriormente, já que são utilizados algoritmos que calculam distâncias entre os atributos, como, por exemplo, IBK, que é um método que tende a dar mais importância para os atributos que possuem um intervalo maior de valores (BATISTA, 2003). Todavia, a normalização não foi eficaz para todos os conjuntos de dados. Os resultados dos dados normalizados são apresentados no Apêndice B.

Com o propósito de comparar o desempenho dos próximos experimentos empregando léxicos, serão utilizados como referência (*baseline*) os melhores resultados de correlação apresentados na Tabela 5.10, para cada traço de personalidade, tendo em vista que, conforme observado na literatura (Seção 2.2), o método TF-IDF é amplamente aplicado no reconhecimento de personalidade por meio de texto.

A próxima seção apresenta os resultados de reconhecimento de personalidade a partir de texto utilizando léxico LIWC.

5.3. Reconhecimento de Traços com LIWC

Nesta seção, os experimentos foram realizados com o léxico LIWC, com o objetivo de validar o reconhecimento de personalidade por meio de texto, de maneira implícita. As categorias do léxico LIWC e métricas utilizadas neste experimento são apresentadas na Seção 3.4.1.

Para avaliar o desempenho dos algoritmos foi utilizado o método de Validação Cruzada com k partes igual a 10. A seguir, são apresentados os resultados de correlação de Pearson com o atributo-meta (traço de personalidade), utilizando como atributos nos conjuntos de testes as categorias do LIWC (conforme descrito na seção 3.4.1). Os valores *baseline* são os melhores resultados da representação vetorial do método TF-IDF, apresentados na seção anterior.

Observa-se que com o léxico LIWC os traços Neuroticismo, Realização e Abertura, obtiveram coeficientes de correlação superiores aos da *baseline*, conforme destaque na Tabela 5.11. Para os demais traços, obtiveram-se resultados próximos aos da *baseline*. Entretanto, a correlação se apresenta no intervalo interpretado como fraco, com resultados entre 0,11 a 0,30.

Tabela 5.11: Resultados de correlação de Pearson do experimento LIWC para todos os traços (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	0.17*	0.1331*	0.1297*	0.2085*	0.058	0.0226	-0.0097
Neuroticismo	0,162	0.1918*	0.1503*	0.1375*	0.0735	0.1738*	0.0531	0.0192
Realização	0,112	0.1131	0.087	0.0503	0.1514*	0.1376*	0.1043	0.0266
Socialização	0,123	-0.0087	-0.0456	-0.0248	0.0405	-0.113	0.017	0.0805
Abertura	0,185	0.1618*	0.1092	0.1693*	0.2008*	0.2031*	0.1206	0.0921

Destaca-se que os valores de correlação conquistados com os experimentos até então realizados, demonstram resultados superiores ao estudo de inferência de personalidade em língua portuguesa descrito na Seção 2.3. A correlação fraca presente nos experimentos e nos trabalhos de RPT para a língua portuguesa corroboram que obter um modelo de inferência de personalidade não é trivial, especialmente em ambientes virtuais, onde o vocábulo utilizado pelos usuários nas redes sociais, na maioria das vezes é informal.

Além disso, o conjunto de treinamento, composto pelas categorias do LIWC, foi submetido ao seletor de atributos “*Wrapper*” e ao método de pesquisa “*Greedy Stepwise*” (detalhes vide Seção 2.2). Para cada algoritmo obteve-se um subconjunto de atributos que possui maior correlação com o atributo-meta. Esses subconjuntos foram utilizados para criar novas bases de treinamentos, contendo apenas as categorias do LIWC que possuem maior ligação com os traços de personalidade.

A Tabela 5.12 apresenta os resultados do conjunto de treinamento LIWC submetidos ao seletor de atributos. Nota-se que com a utilização desse método os resultados superaram a *baseline* para todos os traços, e o coeficiente de correlação aumentou, chegando a uma escala moderada, com valores entre 0,31 a 0,59. Ainda, a tabela destaca o melhor resultado para cada traço.

O método de seleção *wrapper* com o algoritmo IBK não obteve subconjunto de atributos no processo de seleção de característica. Dessa maneira, não foi possível comparar resultados de tal indutor.

Tabela 5.12: Resultados de correlação de Pearson do experimento LIWC com seletor de atributos para todos os traços (* correlação significativa ao nível de 0,05).

Traços	Algoritmos							
	<i>Baseline</i>	M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (<i>KNN=1</i>)	IBK (<i>KNN=3</i>)
Extroversão	0,264	0.0775	0.2539*	0.2595*	0.2709*	0.1656*	--	--
Neuroticismo	0,162	0.2692*	0.2553*	0.2951*	0.3937*	0.1983*	--	--
Realização	0,112	0.2858*	0.2498*	0.2438*	0.3075*	0.2038*	--	--
Socialização	0,123	0.1007	0.0629	0.1628*	--	0.2046*	--	--
Abertura	0,185	0.284*	0.3261*	0.3743*	0.3479*	0.3451*	--	--

Além da comprovação do modelo por meio do coeficiente de correlação, também foram calculados os valores para os testes da raiz quadrada da média do erro (*Root Mean Squared Error* - RMSE) para o experimento utilizando seletor de atributos. Os resultados RMSE, estão dispostos na Tabela 5.13. Observa-se que quanto menor for o RMSE, mais ajustado é o modelo.

Tabela 5.13: RMSE do experimento LIWC com seletor de atributos para todos os traços.

Traços	Algoritmos					
	<i>Baseline</i>	M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL
Extroversão	28.8175	27.9519	27.0434	26.923	26.923	27.6388
Neuroticismo	26.8816	26.2514	26.1837	25.975	24.9958	26.6054
Realização	31.6259	22.102	22.2713	22.4572	22.1236	22.522
Socialização	25.0656	24.8896	24.8692	24.5738	--	24.3912
Abertura	32.9918	24.0322	23.7345	23.3555	23.9663	23.5665

Para determinar a existência (ou não) de diferença estatística significativa entre os algoritmos, no experimento utilizando o léxico LIWC, utilizou-se o teste de Friedman e *post-hoc* de Nemenyi. O teste de Friedman não apontou diferença significativa entre os algoritmos com 95% de confiança. Dessa maneira, nenhum algoritmo utilizado no experimento com o léxico LIWC obteve relevância de desempenho em relação a outro.

Destaca-se que utilizando o léxico LIWC, com a seleção de atributos, os resultados foram superiores à abordagem do vetor de característica TF-IDF, obtendo correlações na escala moderada.

A próxima seção apresenta experimentos associando o léxico LIWC aos léxicos especializados em termos afetivos, com o intuito de obter melhores resultados, por meio de identificação de emoções que corroborem o método proposto de reconhecimento de personalidade a partir de texto.

5.4. Reconhecendo Traços com LIWC Associados a Léxicos Afetivos

Para avaliar se a identificação de emoções contribui com a averiguação dos traços de personalidade, foram utilizados os léxicos afetivos: *SentiStrength*, *AnewBr*, *SentiLex-PT* e *OpLexicon*. Todos os léxicos afetivos citados foram construídos ou adaptados para a língua portuguesa por outros pesquisadores. Ainda, justifica-se a utilização de vários léxicos com o propósito de comparar quais possuem melhores resultados, podendo os mesmos serem combinados entre si.

Os experimentos com os léxicos afetivos foram utilizados em conjunto com o léxico LIWC, sendo divididos em duas abordagens: a primeira contabiliza a polaridade (positivos, negativos e neutros) dos termos empregados nos textos; e a segunda, submete a lista dos termos representativos no texto aos léxicos, considerando o peso de relevância no texto e o peso dos léxicos para o termo. Essas abordagens e a equação de cálculo são dadas na Seção 3.4.3.

Todos os experimentos desta seção foram submetidos ao grupo de indutores utilizados nos experimentos anteriores. No treinamento foi utilizada Validação Cruzada $k = 10$.

Nas próximas subseções são apresentados os resultados obtidos para cada léxico afetivo, juntamente com os resultados da união dos léxicos.

5.4.1. SentiStrength

O primeiro experimento com o léxico *SentiStrength* teve o objetivo de computar o número de palavras positivas e negativas empregadas no texto, ignorando o peso de variação que existe no léxico (intervalo de 1 a 5). Para informações sobre o método consultar a Seção 3.4.3.

A Tabela 5.14, apresenta as estatísticas dos termos positivos e negativos identificados nos textos com o léxico *SentiStrength*. A contabilização dos termos foi feita por usuário. Cumpre esclarecer que a expressão “mínimo” (Tabela 5.14) refere-se ao menor valor encontrado de termos positivos/negativos na base textual entre todos os usuários. Por sua vez,

o termo “máximo” traz o maior valor encontrado da polaridade entre todos os usuários. Em seguida, o termo “Média por usuário” diz respeito à média da polaridade por usuário. Já a estatística “% Mínimo” reporta o menor percentual de termos positivos/negativos, em relação ao total de palavras do texto redigidas pelo usuário, entre todos os usuários do experimento. Por sua vez, o termo “% Máximo” traz o maior percentual de termos positivos/negativos, em relação ao total de palavras do texto, entre todos os usuários. Por fim, a “Média % por usuário” diz respeito à média do percentual de termos positivos/negativos em relação ao total de palavras empregadas nos textos.

Tabela 5.14: Estatística das polaridades do léxico *SentiStrength* identificados nos textos.

Estatísticas	Positivos	Negativos
Mínimo	16	21
Máximo	4.632	5.986
Média por usuário	379	317
% Mínimo	1.41	1.50
% Máximo	10.97	6.67
Média % por usuário	4.52	4.00

Deste modo, ao término da contabilização, foi acrescentado ao conjunto de treinamento do LIWC, a quantidade de palavras positivas e negativas empregadas nos textos de cada indivíduo. O resultado da correlação de Pearson para cada traço, obtido neste experimento, pode ser visualizado na Tabela 5.15. Ainda, a tabela destaca o melhor resultado para cada traço.

Tabela 5.15: Resultados de correlação de Pearson do experimento da primeira abordagem dos léxicos LIWC e *SentiStrength* (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	0.1717*	0.2025*	0.1065	0.22*	0.0483	0.0173	0.0178
Neuroticismo	0,162	0.1587*	0.1355*	0.1312*	0.0766	0.1735*	0.0272	0.0446
Realização	0,112	0.1161	0.0394	0.0242	0.149*	0.1408*	0.0956	0.0358
Socialização	0,123	-0.0112	-0.0618	-0.0339	0.0463	-0.1163	0.0118	0.1034
Abertura	0,185	0.1539*	0.0693	0.1573*	0.2008*	0.2047*	0.1165	0.0552

Observa-se que o experimento contendo os léxicos LIWC e *SentiStrength*, com a aplicação da contabilização de polaridades, auxiliou sutilmente a melhora do método contendo apenas o LIWC para o reconhecimento de personalidade. A correlação obtida se encontra na escala de valores entre 0,11 a 0,30, correspondendo a uma correlação fraca.

A seguir, a Tabela 5.16 apresenta um comparativo entre os resultados do conjunto de treinamento LIWC e o conjunto de treinamento LIWC com *SentiStrength*, utilizando sumarização e polaridades. A tabela destaca os valores do coeficiente de correlação que foram superiores ao conjunto de treinamento contendo apenas o léxico LIWC.

Tabela 5.16: Comparativo entre os resultados de correlação de Pearson entre LIWC e LIWC com *SentiStrength*, utilizando sumarização de polaridades (* correlação significativa ao nível de 0,05).

Traços	Algoritmos						
	M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0.1717*	0.2025*	0.1065	0.22*	0.0483	0.0173	0.0178
Neuroticismo	0.1587*	0.1355*	0.1312*	0.0766	0.1735*	0.0272	0.0446
Realização	0.1161	0.0394	0.0242	0.149*	0.1408*	0.0956	0.0358
Socialização	-0.0112	-0.0618	-0.0339	0.0463	-0.1163	0.0118	0.1034
Abertura	0.1539*	0.0693	0.1573*	0.2008*	0.2047*	0.1165	0.0552

O experimento avança para a segunda abordagem do léxico *SentiStrength*. O conjunto TF-IDF-1500, contendo 1.500 termos com seus pesos de frequências, foi submetido ao léxico com o intuito de averiguar quais desses termos estão contidos no léxico. Por conseguinte, utilizou-se a frequência do termo juntamente com o peso do léxico (intervalo de 1 a 5) para serem associados ao conjunto de treinamento do léxico LIWC. Esse processo resultou em 190 termos identificados no léxico *SentiStrength*.

A seguir são listadas as palavras encontradas no léxico *SentiStrength*, em que cada termo possui uma escala de emoção classificada em um intervalo de 1 a 5, denotando palavras positivas, -5 a -1 para palavras negativas.

Tabela 5.17: Lista de termos TF-IDF identificados no léxico *SentiStrength*.

Termos
abraço, abraços, absurdo, acaso, agradável, alegria, ama, amado, amor, apenas, apesar, apoio, assim, bater, bebe, beijo, bem, boa, bom, bonito, breve, briga, brilho, brincadeira, calma, cansado, carinho, céu, chama, chato, cheio, chora, chorando, chorar, claro, comemorar, compartilhar, completamente, completo, comum, concurso, confiança, contrario, corrupção, crise, cruz, cuidado, cuidar, culpa, desafio, desespero, diabo, difícil, direito, doce, doença, dor, dores, duvida, emoção, errado, errar, erro, escuro, especial, especialmente, espera, esperança, esperando, estranho, faca, fácil, falsidade, falta, favor, fé, feliz, festa, fez, fiel, fome, forca, forte, fraco, frio, fugir, ganhou, gloria, gosta, gostei, gosto, graça, graças, greve, gritar, guerra, honra, humor, idiota, ilusão, incrível, infelizmente, inferno, inimigo, inteligente, inveja, ira, jogar, lagrimas, lixo, louco, loucura, luta, lutar, mal, mano, mar, matar, medo, melhor, melhorar, mentira, mentiras, merda, morrer, morte, namorada, obrigado, odeio, ódio, ok, orgulho, paixão, parar, paz, pecado, pena, perde, perder, perdida, perdido, perfeito, piada, pior, porra, preguiça, problema, pronto, puro, querida, querido, raiva, realmente, respeito, responsável, rir, risco, roubar, ruim, sabedoria, sábio, saída, salvar, saudade, sede, segredo, seguro, serio, simples, sincero, sofrimento, sorriso, sorte, sozinho, sucesso, super, tipo, total, totalmente, triste, tristeza, valor, valores, vão, vazio, velho, verdade, verdadeiro, vergonha, vitória

Após a submissão dos termos no léxico e a associação dos mesmos no conjunto de treinamento do LIWC, o referido conjunto foi disposto aos indutores para averiguar se, com tal abordagem, houve melhora no método. A Tabela 5.18 apresenta os resultados obtidos desse experimento e destaca os valores do coeficiente de correlação que foram superiores ao conjunto de treinamento contendo apenas o léxico LIWC.

Tabela 5.18: Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e *SentiStrength* (* correlação significativa ao nível de 0,05).

Traços	<i>Baseline</i>	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (<i>KNN=1</i>)	IBK (<i>KNN=3</i>)
Extroversão	0,264	0.1878*	0.0298	0.1853*	0.1093	0.0002	0.1126	0.1312*
Neuroticismo	0,162	0.1349*	0.0809	0.1781*	-0.104	0.1079	-0.0992	-0.2061
Realização	0,112	0.1736*	0.016	-0.0467	-0.0293	0.1063	-0.1408	-0.0645
Socialização	0,123	-0.0226	-0.0628	-0.0396	0.0709	-0.036	0.0492	0.032
Abertura	0,185	0.1134	0.0823	0.1475*	0.0806	0.1062	0.0222	0.0535

Visando melhores resultados com essa abordagem, uma redução na dimensionalidade se fez necessária. Para reduzir a dimensão e os ruídos na base, foi realizada uma seleção de atributos baseada nos seguintes parâmetros: “*Wrapper*”, como método seletor e “*Greedy Stepwise*” para determinar o tipo de pesquisa que será realizada. Após esse processo, novos testes foram realizados para cada traço. A Tabela 5.19 apresenta os resultados de correlação dos subconjuntos de atributos selecionados.

Tabela 5.19: Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e *SentiStrength* com seleção de atributos (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	0.4282*	0.5052*	0.4139*	0.5644*	0.1425*	0.2807*	0.2373*
Neuroticismo	0,162	0.2183*	0.3623*	0.3446*	0.3539*	0.2674*	0.1627*	0.1056
Realização	0,112	0.4072*	0.3804*	0.4418*	0.4294*	0.2657*	0.2015*	0.2632*
Socialização	0,123	0.1974*	0.2479*	0.3299*	0.1817*	0.2579*	0.1615*	0.2112*
Abertura	0,185	0.4115*	0.3932*	0.4285*	0.4138*	0.2981*	0.1672*	0.2469*

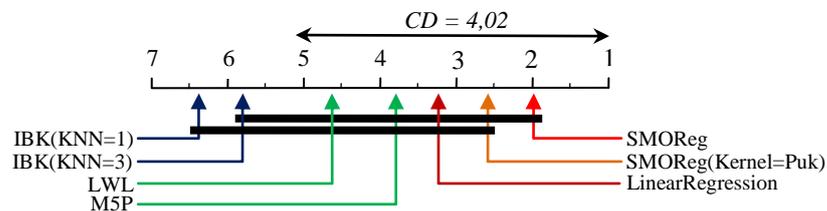


Figura 5.2: Resultado do teste de Nemenyi para o comparativo dos algoritmos do experimento LIWC associado ao léxico *SentiStrength*.

Observa-se que ao aplicar a seleção de atributos, os resultados são mais consistentes, o que mostra que a seleção de características é fundamental para se obter um aumento geral de correlação. Os resultados superaram a *baseline* e para todos os traços de personalidade obteve-se a existência de correlação moderada. Ressalta-se o traço Extroversão, que obteve uma correlação próxima a forte.

Utilizou-se o teste de Friedman para determinar a possibilidade da existência de diferença estatística significativa entre os algoritmos neste conjunto de experimentos. O mesmo apontou que há diferença. Finalmente, utilizou-se o teste *post-hoc* de Nemenyi, cujo resultado é apresentado na Figura 5.2, onde se pode observar que, de maneira indireta, não

existe diferença significativa entre os algoritmos, que justifique a escolha de um ou de outro, exceto IBK ($KNN=1$) que obteve maior *rank* médio.

Uma primeira conclusão importante a se destacar com a utilização dos léxicos LIWC e *SentiStrength* é que a construção de um método capaz de unir léxico psicolinguístico e léxico afetivo corrobora que, para a tarefa de reconhecimento de traços de personalidade por meio de texto, a identificação de emoções contribui com o modelo.

A Subseção 5.4.2 apresentará os experimentos e resultados utilizando o léxico *AnewBr*.

5.4.2. AnewBr

Os experimentos utilizando o léxico *AnewBr*, conforme os demais léxicos afetivos, foram divididos em duas abordagens. A primeira é computar a média ponderada de valência e alerta empregada nos textos e verificar seu desempenho no método. Por sua vez, a segunda abordagem é submeter ao léxico *AnewBr* um conjunto de termos com seus pesos de frequências (TF-IDF), para averiguar quais desses termos estão contidos no léxico e a eles considerar o peso de valência e alerta. Será utilizado o conjunto denominado TF-IDF-1500, contendo 1.500 termos com os valores de frequência. Na Seção 3.4.3 são expostos detalhes sobre essas abordagens.

Para a primeira abordagem, foram adicionados ao conjunto de treinamento LIWC duas novas colunas, valência média e alerta médio, totalizando 70 atributos no conjunto de treinamento. Os indutores foram treinados e testados com esse conjunto de dados. A Tabela 5.20 apresenta os resultados de correlação de Pearson com os traços de personalidade obtidos com tal abordagem.

Tabela 5.20: Resultados de correlação de Pearson do experimento da primeira abordagem dos léxicos LIWC e *AnewBr* (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK ($KNN=1$)	IBK ($KNN=3$)
Extroversão	0,264	0.1834*	0.108	0.1344*	0.209*	0.0104	0.0463	-0.0128
Neuroticismo	0,162	0.1617*	0.1467*	0.1124	0.0748	0.1731*	0.0231	0.0249
Realização	0,112	0.0966	0.0579	0.0613	0.1389*	0.1482*	0.0401	0.0276
Socialização	0,123	0.0101	-0.0541	-0.0327	0.0534	-0.0977	0.0474	0.064
Abertura	0,185	0.168*	0.0882	0.169*	0.2128*	0.1947*	0.1271*	0.0922

Comparando o desempenho entre *AnewBr* e *SentiStrength*, utilizando a primeira abordagem de experimento, é possível observar que ambos os léxicos obtiveram o mesmo nível de correlação (fraca), valores de coeficientes próximos e a maioria dos traços foram superiores aos valores da *baseline*. Destaca-se que o *AnewBr* possui um conjunto de vocábulos inferior ao *SentiStrength*, ou seja, uma cobertura menor de termos afetivos, entretanto, os léxicos alcançaram desempenhos similares. Exceto, para o reconhecimento do traço Socialização, em que a contagem das polaridades do léxico *SentiStrength* foi significativamente superior ao do *AnewBr*.

No próximo experimento, utilizou-se o léxico *AnewBr*, que diz respeito a segunda abordagem, submetendo ao léxico um conjunto de termos contendo seu peso de frequências (TF-IDF). A saber, foram identificados 244 termos no léxico *AnewBr* resultando uma cobertura maior em comparação com o léxico *SentiStrength*, o que significa que o léxico *AnewBr* possui termos com mais relevância para o texto. A Tabela 5.21 apresenta os termos que possuem mais representatividade nos textos identificados no léxico *AnewBr*.

Tabela 5.21: Lista de termos TF-IDF identificados no léxico *AnewBr*.

Termos
abençoado, abraçar, abraço, absurdo, acaso, acordo, agradável, agua, alegre, alegria, amado, amigo, amor, aniversario, anjo, aprender, ar, arte, arvore, azul, banho, bebe, beijo, beleza, belo, bolo, bom, bonito, branco, cabelo, cachorro, calor, cama, campeão, campo, canção, cansado, capaz, carro, casa, casal, casamento, céu, chocolate, chuva, cidade, cinema, comer, computador, confiança, conhecimento, controle, cor, coração, coragem, corpo, costa, criança, crise, desafio, desculpa, desejo, deus, diabo, digno, dinheiro, doce, doença, dor, educação, emprego, encontro, erro, esconder, escuro, espaco, esperanca, espirito, esposa, estranho, estrela, evento, faca, facil, familia, fase, favor, feliz, feriado, ferias, festa, flor, fogo, forca, forte, frio, gato, gloria, gosto, guerra, historia, homem, honra, hospital, humor, ideia, idiota, igreja, inferno, inteligente, irmão, janela, janta, jardim, jogo, justica, lar, leite, letra, liberdade, liga, lindo, livre, livro, lixo, louco, luta, luz, mãe, mal, maneira, mão, maquina, massa, medo, melhorar, memoria, menina, menino, mente, mentira, mercado, mês, mesa, milagre, momento, moral, morte, morto, mulher, mundo, musica, namorada, nascimento, natal, natural, natureza, nome, noticia, novo, ódio, ônibus, opção, opinião, orgulho, ouro, pai, pais, paixão, papel, parte, passagem, paz, pé, pecado, peito, pensamento, perdido, pessoa, piada, pizza, plano, poesia, porta, povo, praia, prazer, presente, preto, primo, problema, processo, professor, prova, qualidade, querido, raiva, rápido, razão, rei, remédio, respeito, resposta, reunião, rico, rio, risada, rocha, rosto, roupa, rua, sábio, salvar, santo, saúde, seguro, sentimento, serio, sexo, social, sociedade, sofrimento, sol, solidão, sombra, sonho, sono, sorriso, sozinho, sucesso, tédio, tempo, terra, tio, triste, vencer, verdade, verde, viagem, vida, vinho, visão, vitória, vivo

Após a submissão dos termos no léxico, o conjunto de treinamento do léxico LIWC associado aos termos identificados no léxico e seus pesos (Frequência e Peso afetivo), foi disposto aos indutores. A Tabela 5.22 apresenta os resultados obtidos desse experimento e destaca os valores do coeficiente de correlação que foram superiores ao conjunto de treinamento contendo apenas o léxico LIWC.

Tabela 5.22: Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e *AnewBr* (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	0.1576*	0.006	0.0337	0.0855	0.0123	0.1631*	0.1714*
Neuroticismo	0,162	0.1393*	0.0189	-0.0172	-0.1108	0.1383*	0.0467	-0.0033
Realização	0,112	0.1823*	-0.0607	-0.0302	-0.0615	0.1169	0.0385	-0.0367
Socialização	0,123	0.0824	-0.0223	0.0005	0.0755	0.005	0.0381	0.0464
Abertura	0,185	0.091	0.0674	0.0544	0.0858	0.1138	0.1522*	0.1436*

Os resultados não foram satisfatórios, nenhuma melhora na correlação foi obtida e escassos casos superaram os valores da *baseline*. A alta dimensionalidade e os ruídos dos enunciados influenciam na obtenção de bons resultados. Dessa forma, o conjunto de treinamento foi submetido à seleção de atributos com o intuito de reduzir a dimensionalidade e os ruídos na base. Os resultados obtidos de correlação dos subconjuntos de atributos selecionados dos léxicos LIWC e *AnewBr*, são apresentados a seguir.

Tabela 5.23: Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e *AnewBr* com seleção de atributos (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	0.1965*	0.5254*	0.3215*	0.5152*	0.3352*	0.0785	0.3738*
Neuroticismo	0,162	0.3612*	0.4539*	0.3993*	0.4778*	0.2132*	0.1627*	0.2181*
Realização	0,112	0.3779*	0.5444*	0.4889*	0.4757*	0.2368*	0.1489*	0.2858*
Socialização	0,123	0.1782*	0.2468*	0.3487*	0.3031*	0.2505*	0.1992*	0.2112*
Abertura	0,185	0.4146*	0.5063*	0.4688*	0.4776*	0.2949*	0.2357*	0.2521*

Com os resultados apresentados, pode-se perceber que, após aplicar a seleção de atributos, os resultados foram mais satisfatórios, o que demonstra que tal técnica é

fundamental para obter-se aumento geral de correlação. Os resultados se sobressaíram quando comparados ao da *baseline*. Ainda, para todos os traços de personalidade, constatou-se a existência de correlação moderada.

Ao comparar as correlações apresentadas com a do léxico *SentiStrength* (subseção anterior), o léxico *AnewBr* obteve melhor resultado na maioria dos traços, quais sejam, Neuroticismo, Realização, Socialização e Abertura. A combinação dos léxicos LIWC e *SentiStrength* só obteve correlação superior, próxima a forte, para o traço Extroversão. Dessa maneira, conclui-se que a combinação dos léxicos LIWC e *AnewBr* possui resultados superiores em relação a combinação LIWC e *SentiStrength*, aplicando seleção de atributos no conjunto de treinamento.

Para identificar o melhor algoritmo ou o melhor grupo de algoritmos, utilizou-se o teste de Friedman, que apontou a existência de diferença estatística significativa entre os algoritmos. Por sua vez, utilizou-se o teste *post-hoc* de Nemenyi para analisar tais diferenças. O resultado é apresentado na Figura 5.3, onde se pode observar que o IBK ($KNN=1$) apresenta maior *rank* médio, o que caracteriza o pior caso.

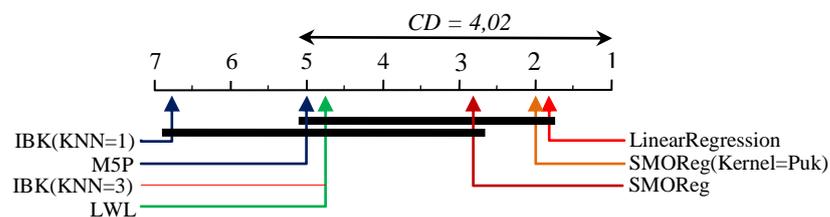


Figura 5.3: Resultado do teste de Nemenyi para o comparativo dos algoritmos do experimento LIWC associado ao léxico *AnewBr*.

5.4.3. SentiLex-PT

O léxico *SentiLex-PT* possui atributos de sentimentos para as polaridades: positivas, negativas e neutras. Dessa forma, o primeiro experimento com o léxico teve o objetivo de computar as polaridades de palavras empregadas no texto. Para informações sobre o método consultar Seção 3.4.3.

Após a contabilização, obtiveram-se as médias de polaridades para cada usuário: 163 por usuário de palavras positivas, 112 para os termos negativos e 60 palavras categorizadas como neutras. A Tabela 5.26, apresenta as estatísticas das polaridades identificadas nos textos com o léxico *SentiLex-PT*. Cumpre esclarecer que a expressão “mínimo”, presente na Tabela

5.24, refere-se ao menor valor encontrado de termos positivos/negativos na base textual entre todos os usuários. Por sua vez, o termo “máximo” traz o maior valor encontrado da polaridade entre todos os usuários. Em seguida, o termo “Média por usuário” diz respeito à média da polaridade por usuário.

Tabela 5.24: Estatística das polaridades do léxico *SentiLex-PT* identificados nos textos.

Estatísticas	Positivos	Negativos	Neutros
Mínimo	8	0	3
Máximo	1.749	1.433	616
Média por usuário	163	112	60

Deste modo, ao término da contabilização, foi acrescentado ao conjunto de treinamento do LIWC a sumarização das polaridades dos termos empregadas nos textos de cada indivíduo, totalizando 72 atributos no conjunto de treinamento, sendo 69 categorias do LIWC e 3 polaridades do léxico. O resultado da correlação de Pearson para cada traço, obtido nesse experimento, pode ser visualizado na Tabela 5.25. Ainda, a tabela destaca o melhor resultado para cada traço.

Tabela 5.25: Resultados de correlação de Pearson do experimento da primeira abordagem dos léxicos LIWC e *SentiLex-PT* (* correlação significante ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	0.1125	0.1545*	0.1546*	0.2068*	0.0469	0.0238	0.0402
Neuroticismo	0,162	0.1902*	0.1943*	0.1272*	0.0922	0.1732*	0.0517	0.0238
Realização	0,112	0.1333*	0.054	0.0456	0.1542*	0.1288*	0.0566	0.0264
Socialização	0,123	0.0623	-0.02	-0.0222	0.0499	-0.1028	0.0387	0.0625
Abertura	0,185	0.1111	0.1096	0.1533*	0.2022*	0.1806*	0.145*	0.0928

O experimento demonstrou que por meio da contabilização das polaridades dos termos não houve diferenças significativas associadas com as categorias do léxico LIWC, ou seja, somente com a utilização das categorias do LIWC obteve-se o mesmo resultado. Os melhores resultados de correlação se encontram na escala de valores entre 0,11 a 0,30, correspondendo a uma correlação fraca.

Ainda, em comparação com a mesma abordagem utilizada para o léxico *SentiStrength*, o léxico *SentiLex-PT* obteve correlações inferiores para os traços Extroversão,

Socialização e Abertura, para os demais traços, obteve correlações superiores, o que demonstra que associação dos léxicos LIWC e *SentiStrength* possui melhores resultados em geral, utilizando a abordagem de sumarização das polaridades dos termos em relação ao conjunto LIWC e *SentiLex-PT*. Isso pode ser justificado pelo fato do léxico *SentiStrength* ter uma maior cobertura dos termos empregados nos textos, com uma média de contabilização de 379 termos positivos por usuário.

Ao comparar os resultados com o experimento do léxico *AnewBr*, não houve diferenças significativas para mensurar qual dos léxicos obteve melhor resultado associado ao LIWC.

Para os léxicos *SentiStrength* e *AnewBr*, os melhores resultados estão associados à utilização da segunda abordagem dos experimentos, que diz respeito a submissão dos léxicos a um conjunto de termos contendo seu peso de frequências (TF-IDF), e à ponderação dos valores de polaridades dos léxicos.

Dessa forma, para o léxico *SentiLex-PT*, foram realizados os mesmos experimentos. Na segunda abordagem foram identificados 179 termos no léxico *SentiLex-PT*. Uma quantidade próxima de termos catalogados no léxico *SentiStrength*. A Tabela 5.26 apresenta os termos do conjunto TF-IDF-1500 que foram identificados no léxico.

Tabela 5.26: Lista de termos TF-IDF identificados no léxico *SentiLex-PT*.

Termos
abençoadado, aberto, absurdo, aceito, acreditar, agradar, agradável, agradecer, alegre, alegria, alto, amado, amar, amigo, animal, arrumar, bater, beleza, belo, bom, bonito, briga, brilho, brincar, cair, calma, cansado, capaz, caralho, carinho, certo, chamado, chato, chorar, confiança, confiar, conquistar, consciência, coragem, corrupção, cuidado, culpa, curto, dado, desespero, desistir, diferente, difícil, digno, direito, direto, doce, educação, energia, enfrentar, engraçado, errado, errar, erro, esconder, escuro, especial, estranho, excelente, experiência, fácil, falsidade, feio, felicidade, feliz, fiel, fofo, fome, forte, fraco, frio, fugir, ganhar, gloria, gostar, gostoso, guerra, honra, humanidade, humano, idiota, igual, ilusão, impossível, incrível, independente, inimigo, inteligência, inteligente, interessante, inveja, jovem, justiça, justo, legal, leve, liberdade, lindo, livre, louco, loucura, lutar, maldade, manha, maravilhoso, mau, medo, melhor, mentira, merda, morte, morto, natural, necessário, normal, objetivo, obrigado, ódio, ótimo, ouvido, paixão, partido, passado, pecado, peço, perder, perdido, perfeito, piada, pior, pobre, preguiça, presente, problema, profissional, puro, quebrar, quente, querido, raiva, rápido, real, respeito, responsável, rico, roubar, ruim, sabedoria, sábio, salvar, santo, segurança, seguro, serio, simples, sincero, sofrer, sofrimento, solidão, sucesso, tédio, triste, tristeza, usar, vão, vazio, velho, vencer, verdade, verdadeiro, vergonha, vitória, vivo, votar

Após a identificação dos termos no léxico e a atribuição de peso (frequência e peso afetivo) aos termos, estes foram associados ao conjunto de treinamento do léxico LIWC e disposto aos indutores. Os resultados não foram satisfatórios, nenhuma melhora na correlação foi obtida e nenhum caso superou a *baseline*, como pode ser observado na Tabela 5.27.

Tabela 5.27: Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e *SentiLex-PT* (* correlação significativa ao nível de 0,05).

Traços	<i>Baseline</i>	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (<i>KNN=1</i>)	IBK (<i>KNN=3</i>)
Extroversão	0,264	0.1869*	0.0493	0.1889*	0.0735	-0.0115	0.1257*	0.0807
Neuroticismo	0,162	0.1575*	0.0156	0.1085	-0.0563	0.0994	-0.0769	-0.0679
Realização	0,112	0.1291*	-0.0532	0.0903	0.0095	0.0974	-0.0236	-0.0291
Socialização	0,123	-0.0559	-0.0936	-0.0875	0.0614	-0.0967	0.0398	0.096
Abertura	0,185	0.1377*	0.0275	0.0067	0.0322	0.115	-0.0267	-0.0075

Conforme aplicação nos demais léxicos, ao utilizar a seleção de atributos no conjunto de treinamento, houve um aumento significativo no coeficiente de correlação com os traços de personalidade. Assim, foi aplicada a seleção de atributos ao conjunto de treinamento desse experimento. Os resultados de tal procedimento podem ser observados na Tabela 5.28.

Nota-se que os valores de correlação se sobressaíram quando comparados ao da *baseline*. Ainda, visualiza-se que para todos os traços de personalidade obteve-se a existência de correlação moderada, as melhores correlações estão entre 0,41 a 0,49.

Tabela 5.28: Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e *SentiLex-PT* com seleção de atributos. (* correlação significativa ao nível de 0,05).

Traços	<i>Baseline</i>	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (<i>KNN=1</i>)	IBK (<i>KNN=3</i>)
Extroversão	0,264	0.4403*	0.483*	0.4818*	0.3676*	0.1576*	0.0785	0.1811*
Neuroticismo	0,162	0.4787*	0.3746*	0.4311*	0.3937*	0.2526*	0.3167*	0.2511*
Realização	0,112	0.4162*	0.3881*	0.3951*	0.2609*	0.2942*	0.1653*	0.1651*
Socialização	0,123	0.1639*	0.1731*	0.1884*	0.4611*	0.1898*	0.1883*	0.2112*
Abertura	0,185	0.3331*	0.4466*	0.4164*	0.4492*	0.2836*	0.2316*	0.3344*

Ao comparar os resultados obtidos com a seleção de atributos aos demais léxicos que também utilizam esse procedimento, o léxico *SentiLex-PT* superou a correlação do léxico *SentiStrength* para os traços Neuroticismo, Socialização e Abertura à experiência. Por sua vez, o léxico *AnewBr* obteve correlações melhores, chegando a uma escala próxima de forte.

Novamente, o teste de Friedman foi utilizado para determinar a existência de diferença estatística significativa entre os desempenhos dos algoritmos neste experimento. O teste de Friedman apontou que não existe diferença entre os algoritmos que justifique a escolha de um ou de outro.

Na próxima subseção serão explorados os experimentos e resultados utilizando o léxico *OpLexicon* para a tarefa de inferência de personalidade.

5.4.4. OpLexicon

O primeiro experimento do presente léxico tem o objetivo de computar o número de palavras positivas, negativas e neutras empregadas pelos participantes em seus discursos, semelhante aos realizados para os demais léxicos discutidos anteriormente.

Após a contabilização das polaridades foram identificadas as estatísticas de cada usuário, com as seguintes médias: 314 palavras positivas; 190 palavras negativas; e 219 palavras neutras. Foi acrescentada a sumarização das polaridades para cada indivíduo ao conjunto de treinamento do LIWC, totalizando um conjunto de 72 atributos. Esse conjunto foi submetido aos algoritmos de aprendizagem de máquina para realizar o treinamento e teste. O resultado da correlação de Pearson para cada traço, obtido neste experimento, pode ser visualizado na Tabela 5.29.

Tabela 5.29: Resultados de correlação de Pearson do experimento da primeira abordagem dos léxicos LIWC e *OpLexicon* (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	0.1408*	0.1636*	0.1532*	0.2035*	0.0458	0.0176	0.0632
Neuroticismo	0,162	0.1836*	0.1782*	0.1358*	0.0903	0.173*	0.0555	0.0142
Realização	0,112	0.1402*	0.0867	0.0499	0.1601*	0.1309*	0.0838	0.0202
Socialização	0,123	0.0422	-0.0643	-0.0206	0.0545	-0.0729	0.0424	0.0676
Abertura	0,185	0.1085	0.1112	0.1624*	0.2058*	0.1797*	0.1568*	0.1339*

Observa-se que o experimento contendo os léxicos LIWC e *OpLexicon*, com a aplicação da contabilização de polaridades, não obteve resultados significativos para o reconhecimento de personalidade. A correlação obtida se encontrada na escala de valores entre 0,11 a 0,30, correspondendo a uma correlação fraca.

Prosseguindo o experimento com o léxico *OpLexicon*, foram identificados quais termos do conjunto TF-IDF (contendo 1.500 palavras) estão contidos no léxico. Baseado na frequência e na polaridade desses termos, esse conjunto foi associado à base do léxico LIWC. A partir da montagem dessa base de dados, o conjunto serviu de entrada para os algoritmos de aprendizagem de máquina utilizados no experimento e obtiveram-se os resultados demonstrados abaixo.

Tabela 5.30: Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e *OpLexicon* (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	0.1195	0.0883	0.0674	0.0511	0.1241*	0.1159	0.0947
Neuroticismo	0,162	0.0726	-0.2029	-0.2048	-0.1306	0.1097	-0.0356	-0.0233
Realização	0,112	0.1469*	0.0995	0.044	-0.0869	0.1198	0.019	0.0225
Socialização	0,123	-0.0693	0.0068	-0.0011	0.065	-0.1474	0.0258	0.0221
Abertura	0,185	0.1666*	0.088	0.0886	0.0514	0.0609	-0.0527	0.1292*

Os referidos resultados não foram satisfatórios, nenhuma melhora na correlação foi obtida e apenas o traço Realização superou os valores da *baseline*. Isto posto, o experimento avançou para etapa de seleção de atributos com o intuito de reduzir a dimensionalidade e consequentemente obter melhores resultados.

Os valores de correlação obtidos com a aplicação de seleção de atributos podem ser visualizados na Tabela 5.31. Nota-se que os valores de correlação se sobressaíram quando comparados a *baseline*, ainda, para os traços Realização e Abertura, obteve-se uma correlação forte, superior a 0,60. Por sua vez, os demais traços de personalidade alcançaram uma correlação moderada, próxima de forte, com valores superiores a 0,55.

Tabela 5.31: Resultados de correlação de Pearson do experimento da segunda abordagem dos léxicos LIWC e *OpLexicon* com seleção de atributos (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	0.4807*	0.58*	0.5483*	0.3813*	0.3136*	0.1907*	0.343*
Neuroticismo	0,162	0.4002*	0.5647*	0.4268*	0.5749*	0.31*	0.1737*	0.1796*
Realização	0,112	0.4864*	0.6264*	0.4588*	0.5312*	0.2942*	0.3051*	0.2044*
Socialização	0,123	0.4399*	0.4631*	0.4581*	0.5572*	0.2757*	0.1912*	0.2112*
Abertura	0,185	0.3992*	0.6013*	0.4334*	0.5749*	0.2616*	0.1624*	0.4192*

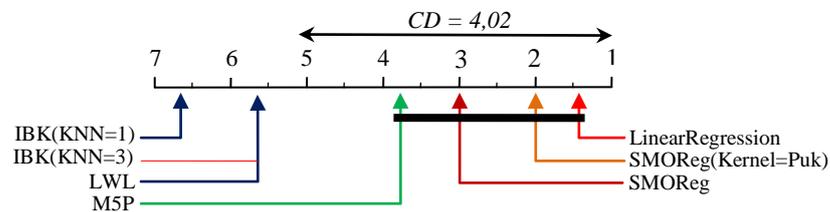


Figura 5.4: Resultado do teste de Nemenyi para o comparativo dos algoritmos do experimento LIWC associado ao léxico *OpLexicon*.

Ao comparar os resultados obtidos a partir da seleção de atributos com demais léxicos, que também utilizam esse procedimento, o método que empregou o *OpLexicon* superou os demais, chegando a uma escala de correlação forte e moderada.

Mais uma vez, o teste de Friedman foi utilizado para determinar a existência de diferença estatística significativa entre os desempenhos dos algoritmos. Esse teste apontou diferença significativa, sendo o resultado do teste de Nemenyi apresentado na Figura 5.4, onde se vê que $\{\text{LinearRegression, SMOReg}(\text{kernel}=\text{puk}), \text{SMOReg, MSP}\} \succ \{\text{LWL, IBK}(\text{KNN}=3), \text{IBK}(\text{KNN}=1)\}$ com 95% de confiança.

Na próxima subseção serão explorados os experimentos utilizando todos os léxicos afetivos juntamente com o LIWC para a tarefa de inferência de personalidade.

5.4.5. Experimentos com a União dos Léxicos Afetivos

Com o intuito de aperfeiçoar o método de reconhecimento de personalidade por meio de texto, optou-se em fazer um experimento com todos os léxicos afetivos citados

anteriormente, com o propósito de abranger uma maior cobertura de termos afetivos empregados no texto, pois cada léxico possui vocábulos diferentes.

Conforme observado nos experimentos individuais de cada léxico, a abordagem que se destaca é aquela que utiliza a frequência do termo no texto (TF-IDF) juntamente com o peso da polaridade no léxico e emprega a seleção de atributos. Dessa maneira, os termos (com seus pesos) de cada léxico afetivo foram identificados no texto e unidos às categorias do LIWC.

Primeiramente, o conjunto de dados foi submetido aos algoritmos para avaliação sem nenhum filtro de característica. A Tabela 5.32 destaca as melhores correlações de Pearson obtidas para esse experimento.

Tabela 5.32: Correlação de Pearson do experimento com todos os léxicos (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	-0.0307	0.1972*	0.179	0.0153	0.1078	0.1074	0.0877
Neuroticismo	0,162	0.0548	-0.0538	-0.0669	-0.1421	0.1423*	0.0161	0.0081
Realização	0,112	0.1004	0.0565	0.0475	-0.1075	0.1214	0.037	-0.0872
Socialização	0,123	0.0186	-0.0294	-0.0197	0.0653	-0.026	0.0549	0.0266
Abertura	0,185	0.0545	0.0998	0.0947	0.0436	0.0378	-0.01	0.0619

Observa-se que os resultados não foram relevantes e não houve correlação superior ao da *baseline*. O aumento expressivo de atributos na base de dados dificultou que o algoritmo de aprendizagem encontrasse um modelo preciso.

O experimento prosseguiu para a etapa de seleção de atributos, obtendo um subconjunto de termos de cada léxico, diminuindo a dimensionalidade do conjunto de dados. Ao submeter o conjunto de dados aos algoritmos para avaliação, com filtro de característica, os resultados foram expressivos. Para os traços de personalidade Extroversão, Neuroticismo e Abertura, a correlação obtida é interpretada como forte, com valores superiores a 0,66. Para os demais traços de personalidade obteve-se a existência de correlação moderada, próxima de forte, com valores superiores a 0,53, conforme demonstrado na Tabela 5.33.

Tabela 5.33: Correlação de Pearson do experimento de seleção de atributos da base de dados com todos os léxicos (* correlação significativa ao nível de 0,05).

Traços	Algoritmos							
	<i>Baseline</i>	M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	0.5723*	0.666*	0.5102*	0.4561*	0.3126*	0.1907*	0.2405*
Neuroticismo	0,162	0.5118*	0.7731*	0.4518*	0.5736*	0.3138*	0.1737*	0.2181*
Realização	0,112	0.5424*	0.5416*	0.5121*	0.5341*	0.2402*	0.3051*	0.2858*
Socialização	0,123	0.533*	0.5176*	0.5869*	0.4518*	0.2572*	0.1992*	0.2112*
Abertura	0,185	0.585*	0.7762*	0.553*	0.6277*	0.2616*	0.1624*	0.4461*

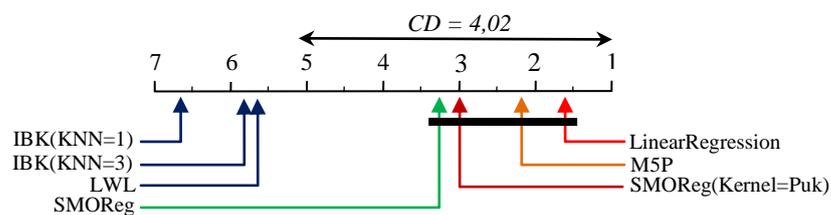


Figura 5.5: Resultado do teste de Nemenyi para o comparativo dos algoritmos do experimento LIWC associado a todos os léxicos afetivos.

Nota-se que os resultados foram satisfatórios, o que corrobora a hipótese de que a união dos léxicos afetivos, combinados entre si, fornece bons resultados para o método de reconhecimento de personalidade por meio de texto. As correlações obtidas com esse experimento são maiores que as correlações dos experimentos individuais de cada léxico, pelo fato de haver maior identificação de termos afetivos empregados nos textos.

Para determinar se existe diferença significativa entre os algoritmos e identificar o melhor algoritmo, aplicaram-se os testes de Friedman e Nemenyi. A Figura 5.5 apresenta o resultado do teste de Nemenyi, onde pode-se ver que $\{\text{LinearRegression, M5P, SMOReg (kernel=puk), SMOReg}\} \succ \{\text{LWL, IBK(KNN=3), IBK(KNN=1)}\}$ com $\alpha = 0.05$.

Devido à obtenção de resultados satisfatórios com a associação dos léxicos, foi levantada a hipótese de que a união desse modelo juntamente com representação vetorial TF-IDF, obterá desempenho significativo no método de reconhecimento de personalidade por meio de texto. Por tal motivo, na próxima seção, serão apresentados os experimentos que verificarão a veracidade de tal hipótese.

5.5. Reconhecendo Traços com LIWC Associados a Léxicos Afetivos e TF-IDF

Com o intuito de averiguar se há melhorias no método de reconhecimento de personalidade, foram associados os conjuntos TF-IDF, apresentados na Seção 5.2, ao léxico LIWC e aos léxicos afetivos.

A cada subconjunto TF-IDF associado aos léxicos, foram realizados testes e avaliação de desempenho dos indutores. Para o conjunto de treinamento, os dados dos léxicos não foram alterados, apenas o conjunto TF-IDF, a ser validado, era alternado. Os conjuntos associados aos léxicos foram: 68, 100, 150, 200, 250, 300, 500, 750 e 1000 termos. Por meio de uma seleção de atributos, todos os conjuntos de dados TF-IDF associados aos léxicos obtiveram uma correlação forte com o atributo-meta que representa o traço de personalidade.

Todavia, serão expostos nesta seção apenas os conjuntos TF-IDF que obtiveram melhores resultados associados aos léxicos, quais sejam: TF-IDF-750 e TF-IDF-1000. Para visualizar todos os resultados de cada conjunto TF-IDF, sugere-se ao leitor consultar o Apêndice B.

O conjunto denominado TF-IDF-750, contém um vetor de característica de 750 termos com seu peso de frequência que corresponde à importância das palavras para o texto. A seguir a Tabela 5.34 ilustra o resultado da associação do conjunto com os léxicos, sem aplicar a seleção de atributos na base de treinamento.

Tabela 5.34: Correlação de Pearson do experimento combinando LIWC, léxicos afetivos e TF-IDF-750, sem seleção de atributos (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK (KNN=1)	IBK (KNN=3)
Extroversão	0,264	0.0915	0.2249*	0.2073*	-0.0242	0.1329*	0.1265*	0.082
Neuroticismo	0,162	0.072	-0.0828	-0.0882	-0.1496	0.0749	-0.0056	0.0102
Realização	0,112	0.0258	0.0164	0.0008	-0.077	0.007	0.0622	-0.0112
Socialização	0,123	0.0199	-0.0473	-0.0418	0.0032	0.0622	0.0205	0.0999
Abertura	0,185	-0.0042	0.0902	0.1064	0.0424	0.0932	0.0835	0.1085

Tendo em vista o número de atributos que coopera para a alta dimensionalidade, observa-se que o experimento que não utilizou seleção de atributos, não reflete resultados relevantes. Todavia, após a etapa de seleção de subconjunto de atributos, os resultados foram

bastante satisfatórios, obtendo correlações superiores a 0,80 para a maioria dos traços, que demonstra uma correlação muito forte entre o atributo-meta e o conjunto de dados.

Para ilustrar os resultados mencionados, a Tabela 5.35 apresenta a correlação de Pearson (r) e a taxa de erro, por meio da raiz quadrada da média do erro (RMSE), dos valores reais e preditos para o atributo-meta (traço), utilizando seleção de atributos.

Tabela 5.35: Correlação de Pearson e RSME do experimento combinando LIWC, léxicos afetivos e TF-IDF-750, com seleção de atributos (* correlação significativa ao nível de 0,05).

Algoritmos	Extroversão		Neuroticismo		Realização		Socialização		Abertura	
	r	RSME								
M5P	0.6079*	22.2058	0.5728*	22.4831	0.6057*	18.5474	0.5478*	20.9314	0.6239*	19.7347
LinearRegression	0.8226*	15.9554	0.8397*	14.8385	0.8145*	13.5989	0.7209*	17.4626	0.8312*	14.1805
SMOReg	0.6037*	22.2807	0.3648*	25.3989	0.5816*	18.9021	0.5225*	21.2896	0.4714*	22.2493
SMOReg (Puk)	0.3876*	25.8756	0.5763*	22.2517	0.4754*	20.492	0.5217*	21.2804	0.5157*	21.6362
LWL	0.2811*	26.709	0.2873*	25.9916	0.2302*	22.6117	0.3308*	23.5892	0.3188*	23.9169
IBK ($KNN=1$)	0.2063*	27.3072	0.2524*	26.8046	0.3025*	22.445	0.1934*	24.6004	0.2522*	24.9571
IBK ($KNN=3$)	0.2919*	27.1373	0.213*	26.5608	0.2832*	22.3224	0.4302*	22.9833	0.445*	22.8767

O segundo conjunto que obteve resultados relevantes é composto por 1.000 termos, contendo o peso de frequência que corresponde a importância de uma palavra para um texto. A seguir a Tabelas 5.36 apresenta os resultados obtidos com esse conjunto de termo, sem a utilização de seleção de atributos no conjunto de treinamento.

Tabela 5.36: Correlação de Pearson do experimento combinando LIWC, léxicos afetivos e TF-IDF-1000, sem seleção de atributos (* correlação significativa ao nível de 0,05).

Traços	Baseline	Algoritmos						
		M5P	Linear Regression	SMOReg	SMOReg (Puk)	LWL	IBK ($KNN=1$)	IBK ($KNN=3$)
Extroversão	0,264	0.1273*	0.2388*	0.2204*	-0.035	0.1393*	0.1433*	0.1836*
Neuroticismo	0,162	0.0864	-0.0687	-0.0734	-0.148	0.0853	-0.0121	0.0065
Realização	0,112	0.0431	0.0457	0.0326	-0.0789	0.0269	0.0604	0.04
Socialização	0,123	0.0238	-0.0403	-0.0325	-0.0013	0.0545	0.0376	0.0615
Abertura	0,185	0.0156	0.1371*	0.1405*	0.0364	0.0724	0.0563	0.1028

Posteriormente, a Tabela 5.37 mostra os resultados de correlação de Pearson (r) e a taxa de erro, por meio da raiz quadrada da média do erro (RMSE), dos valores reais e preditos

para o atributo-meta (traço) do conjunto contendo a combinação do léxico LIWC, léxicos afetivos e TF-IDF-1000, utilizando seleção de atributos.

Tabela 5.37: Correlação de Pearson e RSME do experimento combinando LIWC, léxicos afetivos e TF-IDF-1000, com seleção de atributos (* correlação significativa ao nível de 0,05).

Algoritmos	Extroversão		Neuroticismo		Realização		Socialização		Abertura	
	<i>r</i>	RSME								
M5P	0.653*	21.2278	0.6107*	21.6086	0.5752*	19.0992	0.6055*	19.9523	0.5924*	20.4423
LinearRegression	0.7649*	18.0534	0.8476*	14.6301	0.7742*	14.8107	0.7203*	17.4919	0.8967*	11.2354
SMOReg	0.6104*	22.0739	0.4713*	24.0254	0.5354*	19.6332	0.3917*	22.9733	0.5172*	21.6397
SMOReg (Puk)	0.6735*	21.017	0.5766*	22.2468	0.5488*	19.5197	0.5933*	20.0791	0.6354*	19.6829
LWL	0.3061*	26.4982	0.1646*	26.8513	0.394*	21.4332	0.3419*	23.48	0.3188*	23.9169
IBK (<i>KNN=1</i>)	0.2063*	27.3072	0.3017*	26.1895	0.3065*	22.4049	0.1934*	24.6004	0.3125*	24.1581
IBK (<i>KNN=3</i>)	0.459*	25.3665	0.213*	26.5608	0.417*	21.8895	0.3261*	24.135	0.4405*	22.7922

Observa-se, novamente, que o experimento que não utilizou a seleção de atributos, teve poucos indícios de melhoras de correlação para os traços Extroversão e Abertura, entretanto, esses resultados não refletem relevância. Contudo, após a etapa de seleção de subconjunto de atributos, os resultados foram bastante satisfatórios, obtendo correlações superiores a 0,72 para todos os traços. Sendo que os traços Neuroticismo e Abertura obtiveram uma correlação muito forte.

Destaca-se que os resultados obtidos pelos experimentos reportados nessa seção foram superiores as demais abordagens vistas nas anteriores. Demonstra-se que as combinações das métricas dos léxicos, juntamente com o método TF-IDF, torna o modelo de reconhecimento de personalidade a partir de texto extremamente eficiente, atingindo correlações excedentes a 0,80, categorizando uma correlação muito forte.

A próxima seção apresenta a análise dos resultados, bem como um comparativo dos experimentos apresentados nas seções anteriores.

5.6. Análise dos Resultados

Tendo em vista os resultados apresentados nas seções anteriores, nesta seção será realizada a análise de cada experimento, a comparação das abordagens utilizadas, bem como a demonstração dos melhores resultados para a tarefa de reconhecimento de personalidade por meio de texto.

O primeiro experimento validou o reconhecimento de personalidade por meio de texto, com a utilização do método TF-IDF. Com os experimentos realizados foi possível avaliar uma correlação fraca, obtendo resultados entre 0,11 a 0,30, para os conjuntos de termos TF-IDF. O traço mais bem avaliado foi Extroversão, com uma correlação de 0,26, o que denota que a relação de frequência de palavras importantes no texto influencia a predição dos traços de personalidade de seus autores escritos em língua portuguesa. Ainda, o método TF-IDF foi utilizado como referência (*baseline*) para os demais experimentos, considerando que tal abordagem também foi aplicada em outras línguas para o reconhecimento da personalidade a partir de textos.

Na sequência, foram realizados experimentos com o léxico LIWC, e por meio dos resultados obtidos foi possível confirmar que a forma como as pessoas escrevem fornece abertura para a observação de aspectos emocionais como, por exemplo, a personalidade.

Por meio da avaliação das categorias do LIWC, que extrai características linguísticas e psicolinguísticas, tais como, discordância, ansiedade, palavrões e outros, obtiveram-se resultados de correlações para todos os traços superiores a *baseline* (TFIDF), chegando a uma escala moderada, com valores entre 0,31 a 0,59. O traço Neuroticismo foi o mais bem avaliado, com coeficiente de 0,39.

Destaca-se que os valores de correlação conquistados com esse experimento, demonstraram resultados superiores quando comparados aos estudos apresentados na literatura para a tarefa de inferência de personalidade em língua portuguesa.

Com o intuito de aperfeiçoar o método de inferência de personalidade, por meio de texto, foram realizados experimentos associando o léxico LIWC a léxicos afetivos, que contém vastos vocábulos rotulados que expressão emoções (positivas, negativas e neutras). Dessa maneira, foram utilizados quatro léxicos afetivos da língua portuguesa, sendo eles: *SentiStrength*, *AnewBr*, *Sentilex-PT* e *OpLexicon*. Todos os léxicos foram validados por meio de duas abordagens para mensurar as emoções empregadas no texto.

A primeira abordagem diz respeito à sumarização de palavras afetivas identificadas no texto. A Tabela 5.38 apresenta os melhores resultados de cada léxico utilizando essa abordagem. Observa-se que os resultados não foram satisfatórios, obtendo-se uma correlação fraca para todos os traços e, ainda, para o traço Extroversão nenhum resultado ultrapassou a *baseline*.

Tabela 5.38: Comparação dos resultados de correlação de Pearson entre os léxicos afetivos que utilizaram a sumarização de polaridades.

Traços	LIWC + <i>SentiStrength</i>	LIWC + <i>AnewBr</i>	LIWC + <i>SentiLex-PT</i>	LIWC + <i>OpLexicon</i>
Extroversão	0.22	0.209	0.2068	0.2035
Neuroticismo	0.1735	0.1731	0.1943	0.1836
Realização	0.149	0.1482	0.1542	0.1601
Socialização	0.1034	0.064	0.0625	0.0676
Abertura	0.2047	0.2128	0.2022	0.2058

Por sua vez, na segunda abordagem se obteve melhores resultados. Tal experimento utilizou-se do método TF-IDF para identificar palavras relevantes no texto e atribuir peso aos termos, conforme sua polaridade nos léxicos afetivos. A Tabela 5.39 apresenta os melhores resultados de correlação obtidos para cada léxico afetivo.

Observa-se que com tais experimentos é possível constatar que o léxico *OpLexicon* obteve desempenho superior aos demais, quando associado ao léxico LIWC para a tarefa de reconhecimento de personalidade por meio de texto. Com essa associação de léxicos (LIWC+*OpLexicon*) foram alcançados bons resultados, chegando a uma correlação forte para os traços Realização e Abertura, e uma correlação moderada, próxima de forte, para os demais traços com valores superiores a 0,55.

Tabela 5.39: Comparação dos resultados de correlação de Pearson entre os léxicos afetivos que utilizaram peso de frequência (TF-IDF).

Traços	LIWC + <i>SentiStrength</i>	LIWC + <i>AnewBr</i>	LIWC + <i>SentiLex-PT</i>	LIWC + <i>OpLexicon</i>
Extroversão	0.5644	0.5254	0.483	0.58
Neuroticismo	0.3623	0.4778	0.4787	0.5749
Realização	0.4418	0.5444	0.4162	0.6264
Socialização	0.3299	0.3487	0.4611	0.5572
Abertura	0.4285	0.5063	0.4466	0.6013

Os léxicos afetivos foram combinados entre si e associados ao léxico LIWC. As correlações obtidas com esse experimento são maiores que as correlações dos experimentos individuais de cada léxico, pelo fato que houve maior identificação de termos afetivos empregados nos textos. Para os traços de personalidade Extroversão, Neuroticismo e Abertura

obteve-se uma correlação forte, com valores superiores a 0,66. Para os demais traços observa-se uma correlação moderada, próxima de forte, com valores superiores a 0,53.

Destaca-se que para obter resultados mais promissores a aplicação de seleção de atributos no conjunto de treinamento é fundamental, pois por meio dessa técnica as correlações do experimento foram satisfatórias.

Ressalta-se que o experimento utilizando a combinação do léxico LIWC, léxicos afetivos e o método TF-IDF alcançou resultados de correlação superiores a 0,81 para todos os traços, exceto Socialização que obteve 0,72, o que demonstra uma correlação muito forte entre o atributo-meta e esse conjunto de dados.

Entre os algoritmos de aprendizagem de máquina utilizados para a tarefa de regressão, o *LinearRegression* obteve menor *rank* médio no teste de Friedman, na maioria dos experimentos, o que destaca-o entre os demais. Isso pode ser justificado pelo fato da regressão linear calcular os coeficientes de maneira a minimizar a diferença quadrática entre a saída real (atributo-meta rotulado), mapeando o conjunto de dados em um linear separável para cada traço de personalidade.

O algoritmo com menor desempenho, maior *rank* médio no teste de Friedman, foi o IBK(KNN), que é sensível ao parâmetro k , o que requer alto custo de avaliação, além de ser tendencioso a *overfitting* quando assume que todos os exemplos mais similares (conjunto de treinamento) encontrados são equivalentemente relevantes (UYSAL; GÜVENIR, 1999), por esse motivo, a precisão da predição do modelo pode ser deteriorada.

O desempenho do método desenvolvido nessa pesquisa sobressaiu a todos os estudos da área de RPT para a língua portuguesa, e ainda, a estudos em outras línguas que utilizaram léxicos mais completos, como por exemplo, os trabalhos de (QUERCIA et al., 2012; MAKOVIKJ et al., 2013; BACHRACH et al., 2012; CELLI; ZAGA, 2013; MAIRESSE, 2007). O trabalho de (MAIRESSE, 2007) obteve como maior correlação o coeficiente 0,33, para o traço Abertura à experiência, utilizando textos em língua inglesa e métodos de aprendizagem de máquina (regressão).

5.7. Considerações Finais

Neste capítulo foram apresentados experimentos realizados com o método de RPT para a língua portuguesa. Nessa avaliação, foram utilizados o léxico LIWC, léxicos afetivos e

a representação de texto por meio do método TF-IDF. Ainda, foram combinados léxicos com o intuito de validar métricas para o método de reconhecimento.

Resultados permitem afirmar que os métodos propostos são superiores, quando comparados a outras abordagens de reconhecimento de personalidade a partir de texto para o português brasileiro. De todos os experimentos realizados, destaca-se como melhor métrica a associação dos léxicos com TF-IDF.

O próximo capítulo apresentará as conclusões obtidas deste trabalho, assim como futuros trabalhos.

Capítulo 6

Conclusão e Trabalhos Futuros

O estudo descrito no presente documento propõe um método automático para detecção de personalidade a partir de textos para o português do Brasil, diminuindo a lacuna encontrada no estado da arte para tal idioma. O objetivo foi quantificar os cinco grandes fatores que descrevem a personalidade de um indivíduo (*BigFive*).

Para a realização dessa proposta, foram efetuadas pesquisas acerca das extrações de personalidade por meio de textos, tanto para a língua inglesa, referência na área, quanto para a língua portuguesa, visando à aquisição de bom embasamento teórico, essencial para a construção do método proposto. A partir das pesquisas realizadas, foi observada a lacuna de métodos para a inferência em língua portuguesa. Dessa maneira, uma abordagem linguística com grande quantidade de características textuais, pode ser útil para alcançar um maior desempenho no reconhecimento da personalidade. Além disso, a literatura demonstra que por meio do acréscimo de um léxico afetivo é possível melhorar o modelo de inferência.

Para confirmar essas hipóteses, aplicou-se aos usuários do experimento o inventário NEO-IPIP 120, que permite quantificar os cinco grandes fatores que descrevem a personalidade de um indivíduo (*BigFive*). Dos mesmos usuários, foram solicitadas credenciais da rede social *Facebook* para que informações fossem coletadas automaticamente por um artefato de *software*. Essas informações foram processadas por meio de algoritmos de mineração de textos para a geração de métodos de inferência de personalidade. O método usa exclusivamente uma abordagem de Aprendizagem de Máquina supervisionada para classificar os textos.

Concluída a implementação do método, iniciou-se a etapa de testes e experimentos para validação do trabalho elaborado. Neste contexto, realizou-se diversos experimentos

combinando os léxicos com o intuito de aperfeiçoar a indução dos traços de personalidade. Tais associações contribuíram para que o modelo obtivesse valores de saída fortemente correlacionados à personalidade humana. O método atingiu uma correlação categorizada muito forte, com valores superiores a 0,80 para a maioria dos traços.

Dessa maneira, o método desenvolvido nessa pesquisa sobressaiu a todos os estudos da área de RPT para a língua portuguesa, possuindo os melhores resultados até o momento. Ainda, o método é capaz de superar estudos em outras línguas como, por exemplo, o inglês, que se utiliza de léxicos mais completos.

Como principais contribuições desta dissertação destacam-se: (i) definição da abordagem de reconhecimento de personalidade por meio de texto em língua portuguesa e estrangeira, (ii) desenvolvimento de métodos via léxicos e TF-IDF para o problema de pesquisa de reconhecimento de traços por meio de textos, (iii) definição e confirmação da hipótese de que a identificação de emoções favorece a inferência de traços a partir de texto, (iv) comparação de léxicos afetivos para a tarefa de inferência de personalidade, (v) experimentação do método mostrando a supremacia em relação aos métodos já existentes de reconhecimento de traço para a língua portuguesa, (vi) resultados conclusivos de que combinações de léxicos e métodos de representação de texto levam a resultados significativamente superiores comparados a utilização individual dos léxicos, (vii) construção de uma base textual redigida em língua portuguesa por usuários da rede social *Facebook*, e (viii) avaliação dos algoritmos de AM para o método proposto.

Pode-se ressaltar como trabalhos futuros a adaptação do método desenvolvido para outras línguas, como o inglês. Ainda, a mensuração da personalidade dos usuários de outras redes sociais, como por exemplo, o *Twitter*, podendo tais resultados ser comparados com os obtidos da base textual do *Facebook*. Essa comparação poderá ter o intuito de identificar se os textos curtos do *Twitter* influenciam na tarefa de inferência de personalidade para a língua portuguesa.

Além disso, no que concerne ao método, pesquisas podem ser realizadas para o entendimento da dependência em relação a diferentes abordagens de seleção de características. Também podem ser objetos de estudos as características selecionadas, bem como a relevância dos léxicos utilizados. Ainda, abre-se um horizonte de pesquisa quanto à utilização de técnicas de classificação multi-rótulo e verificação de suas performances.

Por fim, até o momento não foi identificada na língua portuguesa nenhuma ferramenta de mercado que atue no reconhecimento de personalidade em redes sociais, que auxilie a personalização de produtos e sistemas de recomendação baseado em personalidade. Logo, o desenvolvimento deste produto também é uma possibilidade de trabalho futuro.

Referências Bibliográficas

AHA, D. W.; BANKERT, R. L. *A comparative evaluation of sequential feature selection algorithms*. In: Learning from Data, Springer New York, p. 199-206, 1996.

AHA, D. W.; KIBLER, D.; ALBERT, M. K. *Instance-based learning algorithms*. Machine learning, p. 37-66, 1991.

ALAM, F.; STEPANOV, E. A.; RICCARDI, G. *Personality Traits Recognition on Social Network-Facebook*, WCPR (ICWSM-13), Cambridge, MA, USA, 2013.

ALLPORT, F. H.; ALLPORT, G. W. *Personality Traits: Their Classification And Measurement*. In: Journal Of Abnormal And Social Psychology, p. 6–40, 1921.

ALLPORT, G. *Psicologia de la personalidad*. (Personality: A Psychological Interpretation), Argentina, 1961.

ALMEIDA, K. M. Avaliação de personalidade em transtorno afetivo bipolar por meio do estudo de pares de irmãos. Tese (Faculdade de Medicina da Universidade de São Paulo, para obtenção de título de Doutor em Ciência) – USP, São Paulo, 2010.

APPOLINÁRIO, F.. Metodologia da Ciência: filosofia e prática da pesquisa. 1 ed. São Paulo: Editora Thomson, 2006.

ARGAMON, S.; DHAWLE, S.; KOPPEL, M.; PENNEBAKER, J. Lexical predictors of personality type. 2005.

ARROJU, M.; HASSAN, A.; FARNADI, G. *Age, Gender and Personality Recognition using Tweets in a Multilingual Setting*, 2015.

ARYA, A.; RAGINI, S.; KUMAR, H.; ABINAYA, G. *A text analysis based seamless framework for predicting human personality traits from social networking sites*. In: International Journal of Information Technology and Computer Science (IJITCS), 2012.

ATKESON, C. G.; MOORE, A. W.; SCHAAL, S. *Locally weighted learning for control*. In: Lazy learning, Springer Netherlands, p. 75-113, 1997.

BACHRACH, Y.; KOSINSKI, M.; GRAEPEL, T.; KOHLI, P.; STILLWELL, D. *Personality and patterns of Facebook usage*. In: proceedings of the 3rd annual ACM web science conference. ACM, p. 24-32. 2012.

BAI, S.; ZHU, T.; CHENG, L. *Big-Five Personality Prediction Based on User Behaviors at Social Network Sites*. In: eprint arXiv:1204.4809, 2012. Disponível em: <<http://arxiv.org/abs/1204.4809v1>>

BAISE, M. Avaliação dos traços de personalidade em pacientes com anorexia nervosa, segundo o Inventário de Temperamento e Caráter de Cloninger. Universidade de São Paulo, 2008. Disponível em: <http://www.teses.usp.br/teses/disponiveis/47/47135/tde-20032009-124600/publico/baise_me.pdf>

BALAGE FILHO, P. P.; PARDO, T. A.; ALUISIO, S. M. *An evaluation of the brazilian portuguese liwc dictionary for sentiment analysis*. In: 9th Brazilian Symposium in Information and Human Language Technology, Fortaleza, Ceara, p. 215-219, 2013.

BARRICK, M. R.; MOUNT, M. K.; JUDGE, T. A. *Personality and performance at the beginning of the new millennium: What do we know and where do we go next?*. In International Journal of Selection and Assessment, v. 9, n. 1-2, p. 9-30, 2001.

BATISTA, G. E. A. P. A. Pré-processamento de Dados em Aprendizado de Máquina Supervisionado. Tese (Instituto de Ciência Matemática e de Computação da Universidade de São Paulo, para obtenção de título de Doutor em Ciência) – USP, São Paulo, 2003.

BENCHERIF, M. A.; ALSULAIMAN, M.; MUHAMMAD, G.; ALI, Z.; MAHMOOD, A.; FAISAL, M. *Gender Effect in Trait Recognition*. In: Proceedings of the World Congress on Engineering and Computer Science, v. 1, 2012.

BERGER, K. S. *The Developing Person Through The Life Span*. 6^a Ed. Worth Publishers, 2003.

BIEL, J.;GATICA-PEREZ, D. *The Youtube lens: Crowdsourced personality impression and audiovisual of vlogs*. In IEEE Transactions on Multimedia, vol. 15, no. 1, pp. 41–55, 2012.

BOUCHET, F.; SANSONNET, J. P. *Classification of wordnet personality adjectives in the NEO PI-R taxonomy*. In: Fourth Workshop on Animated Conversational Agents, p. 83-90, 2010.

BRADLEY, M. M.; LANG, P. J. *Affective Norms for English Words (ANEW): instruction manual and affective ratings*. 1999.

BRADLEY, M. M.; LANG, P. J. *Measuring emotion: The self-assessment manikin and the semantic differential*. In: Journal of Behavior Therapy and Experimental Psychiatry, p. 49–59, 1994.

BUSSAB, W. O; MORETTIN, A. *Estatística Básica*. 3 ed. São Paulo: Atual, 1986.

CAMBRIA, E.; LIVINGSTONE, A.; HUSSAIN, A. *The hourglass of emotions*. In: Cognitive Behavioural Systems, p. 144–157, 2012.

CAMBRIA, E.; OLSHER, D.; RAJAGOPAL, D. *SenticNet 3: A common and common-sense knowledge base for cognition-driven sentiment analysis*. In: AAI, pp. 1515-1521, Quebec City, 2014.

CAMBRIA, E.; SPEER, R.; HAVASI, C.; HUSSAIN, A. *SenticNet: A Publicly Available Semantic Resource for Opinion Mining*. In: AAAI fall symposium: commonsense knowledge, v. 10, p. 02, 2010.

CAPRARA, G. V.; BARBARANELLI, C.; BORGOGNI, L.; PERUGINI, M. *The “Big Five Questionnaire”*: A new questionnaire to assess the Five Factor Model, In *Personality and Individual Differences*, vol. 3, p. 281-288, 1993.

CARVALHO, P.; SILVA, M. J. SentiLex-PT: Principais características e potencialidades. In: *Oslo Studies in Language*, 2015.

CAVNAR, W. B.; TRENKLE, J. M. *N-gram-based text categorization*, In: *Ann Arbor MI*, vol. 2, p. 161-175, 1994.

CELLI, F. *Adaptive Personality Recognition from Text*, PhD Thesis, University of Trento(Italy), 2012.

CELLI, F. *Mining user personality in twitter*. In: *Language, Interaction and Computation CLIC*, 2011.

CELLI, F.; BRUNI, E.; LEPRI, B. *Automatic Personality and Interaction Style Recognition from Facebook Profile Pictures*. In: *Proceedings of the ACM International Conference on Multimedia*, p. 1101-1104, 2014. ISBN: 978-1-4503-3063-3. Disponível em: <<http://dx.doi.org/10.1145/2647868.2654977>>.

CELLI, F.; ZAGA, C. *Be Conscientious, Express your Sentiment!*. In: In *Proceedings of ESSEM*, in conjunction with aixia 2013, p. 140-147, 2013.

CHIKERSAL, P.; PORIA, S.; CAMBRIA, E. *SeNTU: Sentiment analysis of tweets by combining a rule-based classifier with supervised learning*. In: *Proceedings of the International Workshop on Semantic Evaluation (SemEval)*, 2015.

COLTHEART, M. *The MRC psycholinguistic database*. In Quarterly Journal of Experimental Psychology, 1981.

CORDER, G.; FOREMAN, D. *Nonparametric Statistics for Non-Statisticians: A Stepby-Step Approach*. Wiley, 2011. ISBN 9781118211250. Disponível em: <<http://books.google.com.br/booksid=T3qOqdpSz6YC>>.

DAMÁSIO A. R. O erro de Descartes emoção razão e o cérebro humano. São Paulo, SP: Companhia das Letras, 1996.

DAVIDOV, D.; TSUR, O.; RAPPOPORT, A. *Enhanced sentiment learning using twitter hashtags and smileys*. In: Proceedings of the 23rd International Conference on Computational Linguistics: Posters, p. 241-249, 2010.

DE RAAD, B.; PERUGINI, M. *Big Five Assessment*, Hogrefe Huber, Ashland, USA, 2002.

DENECKE, K. *Using sentiwordnet for multilingual sentiment analysis*. In: Data Engineering Workshop, IEEE 24th International Conference on, p. 507-512, 2008.

DIAS-DA-SILVA, B. C.; DI FELIPPO, A.; NUNES, M. G. V. *The automatic mapping of Princeton WordNet lexical conceptual relations onto the Brazilian Portuguese WordNet database*. In: Proceedings of the 6th International Conference on Language Resources and Evaluation (LREC 2008), Marrocos, 2008.

DIMURO, G. P.; COSTA, A. C. R., GONÇALVES, L. V.; HÜBNER, A. *Centralized regulation of social exchanges between personality-based agents*. In: Coordination, Organizations, Institutions, and Norms in Agent Systems II, Springer Heidelberg, Berlin, Germany, p. 338-355, 2007.

DWECK, C. S. *Can personality be changed? The role of beliefs in personality and change*. In Current Directions in Psychological Science, v. 17, n. 6, p. 391-394, 2008.

EMOSENTICNET. *EmoSenticNet*. Disponível em: <<http://www.gelbukh.com/emosenticnet/>>
Acesso em: 24 abr. 2015.

ESULI, A.; SEBASTIANI, F. *Sentiwordnet: A publicly available lexical resource for opinion mining*. In: Proceedings of LREC, p. 417-422, 2006.

FARNADI, G., SUSHMITA, S., SITARAMAN, G., TON, N., DE COCK, M., & DAVALOS, S. *A Multivariate Regression Approach to Personality Impression Recognition of Vloggers*. In: Proceedings of the 2014 ACM Multi Media on Workshop on Computational Personality Recognition, p. 1-6, 2014.

FARNADI, G.; ZOGHBI, S.; MOENS, M. F.; DE COCK, M. *Recognising personality traits using Facebook status updates*. In: Proceedings of the Workshop on Computational Personality Recognition (WCPR), AAAI Press, p. 14–18, 2013.

FERWERDA, B.; YANG, E.; SCHEDL, M.; TKALCIC, M. *Personality traits predict music taxonomy preferences*. In: Proceedings of the 33rd Annual ACM Conference Extended Abstracts on Human Factors in Computing Systems, p. 2241-2246, 2015.

FRAKES, W. B. *Stemming algorithm*. Information Retrieval: data structures and algorithms, Prentice Hall, Englewood Cliffs, New Jersey, USA, 1992.

FRANK, E.; TRIGG, L.; HOLMES, G.; WITTEN, I. H. *Technical note: Naive Bayes for regression*. Machine Learning, v. 41, n. 1, p. 5-25, 2000.

FREITAS, L. A.; VIEIRA, R. *Ontology based feature level opinion mining for portuguese reviews*. In: Proceedings of the 22nd international conference on World Wide Web companion, p. 367-370, 2013.

FREUND, Y.; IYER, R.; SCHAPIRE, R. E.; SINGER, Y. *An efficient boosting algorithm for combining preferences*. In: Proceedings of the 15th International Conference on Machine Learning, p. 170–178, 1998.

FRIEDMAN, M. *The Use of Ranks to Avoid the Assumption of Normality Implicit in the Analysis of Variance*. Journal of the American Statistical Association, American Statistical Association, v. 32, n. 200, p. 675–701, dez. 1937. ISSN 01621459. Disponível em: <<http://dx.doi.org/10.2307/2279372>>.

GARCIA, D.; SCHWEITZER, F. *Emotions in Product Reviews--Empirics and Models*. In Privacy, Security, Risk and Trust (PASSAT) and 2011 IEEE Third International Conference on Social Computing (SocialCom), p. 483-488, 2011.

GARCIA, L. F. Teorias psicométricas da personalidade. In: Introdução à psicologia das diferenças individuais, p. 219-242, 2006.

GEZICI, G., DEHKHARGHANI, R., YANIKOGLU, B., TAPUCU, D., & SAYGIN, Y. *Su-sentilab: A classification system for sentiment analysis in twitter*. In: Proceedings of the International Workshop on Semantic Evaluation, p. 471-477, 2013.

GILL, A. J.; NOWSON, S.; OBERLANDER, J. *What Are They Blogging About? Personality, Topic and Motivation in Blogs*. In: Third International AAAI Conference on Weblogs and Social Media, 2009.

GOLBECK, J.; ROBLES, C.; EDMONDSON, M.; TURNER, K. *Predicting personality from twitter*. In: International Conference on Social Computing (SocialCom), IEEE Third International Conference, p. 149-156, 2011.

GONÇALVES, P.; ARAÚJO, M.; BENEVENUTO, F.; CHA, M. *Comparing and combining sentiment analysis methods*. In: Proceedings of the first ACM conference on online social networks, p. 27-38, 2013.

GOSLING, S. D.; RENTFROW, P. J.; SWANN J. R., W. B. *A very brief measure of the big-five personality domains*. In: Journal of Research in Personality, p. 504–528, 2003.

GUNTUKU, S. C.; QIU, L.; ROY, S.; LIN, W.; JAKHETIYA, V. *Do Others Perceive You As You Want Them To? Modeling Personality based on Selfies*. In: Proceedings of the 1st International Workshop on Affect & Sentiment in Multimedia, p. 21-26, 2015. Disponível em: <<http://dx.doi.org/10.1145/2813524.2813528>>.

HALL M. A. *Correlation-based Feature Subset Selection for Machine Learning*. Hamilton, New Zealand. 1998.

HALL, M. A.; SMITH, L. A. *Practical feature subset selection for machine learning*. In Proceedings of Australasian Computer Science Conference. Springer, Singapore, p. 181–191, 1998.

HALLIDAY, M. A. K. *Functional grammar*. London: Edward Arnold, 1994.

HENDRIKS, A. A. J. *The construction of the Five-Factor Personality Inventory (FFPI)*. University of Groningen, Holanda, 1997.

HENNA, E. A. D. *Relação entre temperamento, caráter e bem-estar subjetivo: estudo em uma amostra de sujeitos saudáveis*. Tese (Faculdade de Medicina da Universidade de São Paulo, para obtenção de título de Doutor em Ciência) – USP, São Paulo, 2011.

HU, M.; LIU, B. *Mining and summarizing customer reviews*. In: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining, p. 168–177, New York, USA, 2004.

HUTZ, C. S.; NUNES, C. H.; SILVEIRA, A. D.; SERRA, J.; ANTON, M.; WIECZOREK, L. S. *O Desenvolvimento De Marcadores Para A Avaliação Da Personalidade No Modelo Dos Cinco Grandes Fatores*. Psicologia: Reflexão E Crítica, 1998.

IACOBELLI, F.; GILL, A. J.; NOWSON, S.; OBERLANDER, J. *Large scale personality classification of bloggers*. In: Affective Computing and Intelligent Interaction, p. 568-577, 2011.

JOHN, O. P.; DONAHUE, E. M.; KENTLE, R. L. *The “Big Five” Inventory: Versions 4a and 5b*. Tech. Institute of Personality and Social Research, University of California, 1991.

JOHNSON, J. A. *Measuring thirty facets of the Five Factor Model with a 120-item public domain inventory: Development of the IPIP-NEO-120*, In: *Journal of Research in Personality*, p. 78-89, 2014.

JOHNSON, J. A. *Web-based personality assessment*, In 71st Annual Meeting of the Eastern Psychological Association, Baltimore, USA, 2000.

KAJI, N.; KITSUREGAWA, M. *Automatic Construction of Polarity-tagged Corpus from HTML Documents*. n. July, p. 452–459, 2006.

KAJI, N.; KITSUREGAWA, M. *Building Lexicon for Sentiment Analysis from Massive Collection of HTML Documents*. In: *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, v. 43, p. 1075–1083, 2007.

KERMANIDIS, K. L. *Mining Authors' Personality Traits from Modern Greek Spontaneous Text*. In: *4th International Work-shop on Corpora for Research on Emotion Sentiment & Social Signals*, in conjunction with LREC12. 2012.

KHALID, S.; KHALIL, T.; NASREEN, S. *A survey of feature selection and feature extraction techniques in machine learning*. In: *Science and Information Conference (SAI)*, IEEE, p. 372-378, 2014.

KIPPER-SCHULER, K. *VerbNet: A Broad-Coverage, Comprehensive Verb Lexicon*. PhD Thesis, University of Pennsylvania, 2005.

KOHAVI, R.; JOHN, G. *Wrappers for feature subset selection*. *Artificial Intelligence*, Amsterdam, v. 97, n. 1-2, p. 273-324, 1997.

KOSINSKI, M.; STILLWELL, D.; GRAEPEL, T. *Private traits and attributes are predictable from digital records of human behavior*. In: Proceedings of the National Academy of Sciences, 2013.

KRISTENSEN C. H; GOMES C. F. A; JUSTO A. R; VIEIRA K. Normas brasileiras para o Affective Norms for English Words. In: Trends Psychiatry Psychother, 2011.

LEVIN, B. *English Verb Classes and Alternation, A Preliminary Investigation*. In: The University of Chicago Press, 1993.

LIKERT, R. *A technique for the measurement of attitudes*. Archives of psychology, 1932.

LIMA, A. C. E. S.; CASTRO, L. N. *Multi-label Semi-supervised Classification Applied to Personality Prediction in Tweets*. In: Computational Intelligence and 11th Brazilian Congress on Computational Intelligence (BRICS-CCI & CBIC), 2013 BRICS Congress on. IEEE, p.195-203, 2013. Disponível em: <<http://dx.doi.org/10.1109/BRICS-CCI-CBIC.2013.41>>

LITVINOVA, T. A.; SEREDIN, P. V.; LITVINOVA, O. A. *Using Part-of-Speech Sequences Frequencies in a Text to Predict Author Personality: a Corpus Study*. In: Indian Journal of Science and Technology, p. 93-97, 2015.

LIU, B. *Sentiment Analysis and Opinion Mining*. Cambridge University Press, 2012.

LUYCKX, K.; DAELEMANS, W. *Using syntactic features to predict author personality from text*, in: Proc. Digit. Humanities, p. 146–149, 2008.

MAIRESSE, F.; WALKER, M. A.; MEHL, M. R.; MOORE, R. K. *Using linguistic cues for the automatic recognition of personality in conversation and text*. In: Journal of artificial intelligence research, p. 457-500, 2007.

MARKOVIKJ, D.; GIEVSKA, S.; KOSINSKI, M.; STILLWELL, D. *Mining facebook data for predictive personality modeling*. In: Proceedings of the 7th international AAAI conference on Weblogs and Social Media, Boston, MA, USA, 2013.

MATHEWS, G.; DEARY, I. J., WHITEMAN, M. C. *Personality Traits*. Cambridge university press, 2009. ISBN 0521831075.

MAZIERO, E. G.; PARDO, T. A. S.; FELIPPO, A. D.; DIAS-DA-SILVA, B. C. A base de dados lexical e a interface web do TeP 2.0: thesaurus eletrônico para o Português do Brasil. In: Proceedings of (WebMedia '08). ACM, New York, NY, USA, 390-392, 2008.

MCCRAE, R. R.; COSTA, P. T. *A five-factor theory of personality*. *Handbook of personality: Theory and research*, v. 2, p. 139-153, 1999.

MCCRAE, R. R.; JOHN, O. P. *An introduction to the five-factor model and its applications*. *Journal of personality*, v. 60, n. 2, p. 175-215, 1992.

MCDUGALL, W. *Of the words character and personality*. In: *Journal of Personality*, v. 1, n. 1, p. 3-16, 1932.

MEDHAT, W.; HASSAN, A.; KORASHY, H. *Sentiment analysis algorithms and applications: A survey*. *Ain Shams Engineering Journal*, v. 5, n. 4, p. 1093–1113, 2014.

MEHL, M. R.; GOSLING, S. D.; PENNEBAKER, J. W. *Personality in its natural habitat: Manifestations and implicit folk theories of personality in daily life*. In: *Journal of Personality and Social Psychology*, p. 862–877, 2006.

MOFFAT, D. *Personality Parameters and Programs*. In: *Creating Personalities for Synthetic Actors*, (Ed.) Trappl, Springer. 1997.

MOHAMMAD, S.M.; KIRITCHENKO, S. *Using Nuances of Emotion to Identify Personality*. 2012.

MONTOYO, A; MARTINEZ-BARCO, P; BALAHUR, A. *Subjectivity and sentiment analysis: An overview of the current state of the area and envisaged developments*. In: *Journal Decision Support Systems*, v.53, ed.4, p. 675-679, 2012.

MOSTAFA, M. M. *An emotional polarity analysis of consumers' airline service tweets*. In: *Social Network Analysis and Mining*, v. 3, p. 635–649, 2013.

MOURÃO, M.; SAIAS, J. BCLaaS: implementação de uma base de conhecimento linguístico as-a-service, 2013.

NETO, A. F. B. *Uma arquitetura para agentes inteligentes com personalidade e emoção*. Tese de Doutorado. Microsoft Research Sao Paulo. 2010. Disponível em: <www.ime.usp.br/~fcs/LIDET/bressane.pdf>

NIELSEN, F. A. *A new ANEW: Evaluation of a word list for sentiment analysis in microblogs, 2011*.

NORMAN, W. T. *Toward an adequate taxonomy of personality attributes: replicated factor structure in peer nomination personality rating*. In: *Journal of Abnormal and Social Psychology*, p. 574–583, 1963.

NOWSON, S.; OBERLANDER, J. *Identifying more bloggers: Towards large-scale*, in: *Proc. Int. Conf. Weblogs Social Media*, 2007.

NUNES, M. A. S. N. *Computação Afetiva personalizando interfaces, interações e recomendações de produtos, serviços e pessoas em ambientes computacionais*. In: *Projetos e Pesquisas em Ciência da Computação no DCOMP/PROCC/UFS*, v. 1, p. 115-151, 2012.

NUNES, M. A. S. N. *Recommender system based on personality traits*. Tese (Université Montpellier, para obtenção de título de Doutor em Informatica), França, 2008.

NUNES, M. A. S. N.; TELES, F. R.; DE SOUZA, J. G. Inferindo personalidade via tweets. In: *GEINTEC-Gestão, Inovação e Tecnologias* 3.3, p. 045-057. 2013.

NUNES, M. A. S. N.; CERRI, S. A.; BLANC, N. *Towards user psychological profile*. In: Proceedings of the VIII Brazilian Symposium on Human Factors in Computing Systems, Sociedade Brasileira de Computação, p. 196-203, 2008.

OBERLANDER, J.; NOWSON, S. *Whose thumb is it anyway? classifying author personality from weblog text*. In: Proceedings of the 44th Annual Meeting of the Association for Computational Linguistics ACL. 2006.

OPLEXICON. *OpLexicon*. Disponível em: <http://ontolp.inf.pucrs.br/Recursos/downloads-OpLexicon.php> . Acesso em: 25 abr. 2015.

ORTONY, A.; CLORE, G. L.; FOSS, M. A. *The Referential Structure of the Affective Lexicon*, 1987.

PALMER, M. *A Class-Based Verb Lexicon*, Disponível em: <<http://verbs.colorado.edu/~mpalmer/projects/verbnet.html>>, Acesso em: 19 de Abril de 2015.

PANG, B.; LEE, L.; VAITHYANATHAN, S. *Thumbs up?: sentiment classification using machine learning techniques*. In: Conference on Empirical methods in natural language processing - EMNLP '02, v. 10, p. 79–86, 2002.

PASQUALI, L. Técnicas de exame psicológico – TEP, ed. Casa do Psicólogo, São Paulo, 2001.

PASQUALOTTI, P. R.; VIEIRA, R. WordnetAffectBR: uma base lexical de palavras de emoções para a língua portuguesa, 2008.

PENG, K. H.; LIOU, L. H.; CHANG, C. S.; LEE, D. S. *Predicting personality traits of Chinese users based on Facebook wall posts*. In: Wireless and Optical Communication Conference (WOCC), p. 9-14, 2015.

PENNEBAKER, J. W.; CHUNG, C. K.; IRELAND, M.; GONZALES, A.; BOOTH, R. J. *The development and psychometric properties of LIWC2007*. 2007.

PENNEBAKER, J. W.; FRANCIS, M. E.; BOOTH, R. J. *Inquiry and Word Count: LIWC*. 2001.

PENNEBAKER, J. W.; KING, L. A. *Linguistic styles: Language use as an individual difference*. In: Journal of Personality and Social Psychology, p. 1296–1312, 1999.

PERVIN, L. A.; JOHN, O. P. *Personalidade - Teoria e Pesquisa*. Ed. Artmed, Porto Alegre, Brasil, 2003.

PICARD, R. W. *Affective Computing*. MIT Press, Cambridge, MA, USA, 1997.

PICCHI NETTO, O. Um filtro iterativo utilizando árvores de decisão, Dissertação de Mestrado, USP, 2013.

POHJALAINEN, J.; RÄSÄNEN, O.; KADIOGLU, S. *Feature selection methods and their combinations in high-dimensional classification of speaker likability, intelligibility and personality traits*. In: Computer Speech & Language, v. 29, p. 145-171, 2015.

PORIA, S.; GELBUKH, A.; AGARWAL, B.; CAMBRIA, E.; HOWARD, N. *Common sense knowledge based personality recognition from text*. In: Advances in Soft Computing and Its Applications, Springer Heidelberg, Berlin, Germany, p. 484-496, 2013.

PORIA, S.; GELBUKH, A.; CAMBRIA, E.; HUSSAIN, A.; HUANG, G. B. *EmoSenticSpace: A novel framework for affective common-sense reasoning*. In: Knowledge-Based Systems, p108-123, 2014. Disponível em: <<http://dx.doi.org/10.1016/j.knosys.2014.06.011>>.

PORIA, S.; GELBUKH, A.; HUSSAIN, A.; DAS, D.; BANDYOPADHYAY, S. *Enhanced SenticNet with Affective Labels for Concept-based Opinion Mining*, Intelligent Systems, IEEE, v. 28, n. 2, p. 31–38, 2013. Disponível em: <<http://doi:10.1109/MIS.2013.4>>.

PORTO, S. M.; COSTA, W. S.; NUNES, M. A. S. N.; MATOS, L. N. Como a extração de personalidade através do teclado pode beneficiar a personalização na Educação. In: Towards Affective Computing in Education (SBIE-WIE 2011 Workshop), p. 1800-1807, 2011.

QIU, L.; LIN, H.; RAMSAY, J.; YANG, F. *You are what you tweet: Personality expression and perception on twitter*, In: Journal of Research in Personality, vol. 46, n. 6, p. 710–718, 2012.

QUERCIA, D.; KOSINSKI, M.; STILLWELL, D.; CROWCROFT, J. *Our Twitter profiles, our selves: Predicting personality with Twitter*. In: Privacy, Security, Risk and Trust (PASSAT) and 2011 IEEE Third International Conference on Social Computing (SocialCom), 2011 IEEE Third International Conference, p. 180-185, 2011.

QUERCIA, D.; LAMBIOTTE, R.; STILLWELL, D.; KOSINSKI, M.; CROWCROFT, J. *The personality of popular facebook users*. In: Proceedings of the ACM 2012 conference on computer supported cooperative work, p. 955-964, 2012.

REEVES, B.; NASS, C. *The media equation: how people treat computers, television, and new media like real people and places*. Cambridge University Press, New York, NY, USA, 1996.

RIGBY, P. C.; HASSAN, A. E. *What can oss mailing lists tell us? A preliminary psychometric text analysis of the apache developer mailing list*. In: Proceedings of the Fourth International Workshop on Mining Software Repositories, p. 23, 2007.

RILOFF, E. *Little words can make a big difference for text classification*. In: Proceedings of the 18th annual international ACM SIGIR conference on Research and development in information retrieval, p. 130-136, 1995.

ROSHCHINA, A.; CARDIFF, J.; ROSSO, P. *A comparative evaluation of personality estimation algorithms for the twin recommender system*. In: Proceedings of the 3rd international workshop on Search and mining user-generated contents, p. 11-18, 2011.

SALTON, G.; BUCKLEY, C. Term-weighting approaches in automatic text retrieval. *Information Processing & Management*, v. 24, n. 5, p. 513–523, 1988.

SANSONNET, J. P.; BOUCHET, F. *Extraction of agent psychological behaviors from glosses of wordnet personality adjectives*. In: Proc. of the 8th European Workshop on Multi-Agent Systems (EUMAS'10), 2010.

SANTOS, M. T.; FLORES-MENDOZA, C. E. Adaptação do Eysenck Personality Questionnaire Júnior para pré-escolares: versão heterorrelato. *Avaliação Psicológica*, v. 11, n. 2, p. 203-212, 2012.

SCARTON, C. E. *VerbNet.Br: construção semiautomática de um léxico verbal online e independente de domínio para o português do Brasil*. Universidade de São Paulo, 2013.

SCARTON, C. E.; ALUISIO, S. *Towards a cross-linguistic VerbNet-style lexicon for Brazilian Portuguese*. In: Workshop on Creating Cross-language Resources for Disconnected Languages and Styles Workshop Programme, p.11, 2012.

SCHIMIT, M. J.; KIHM, J. A.; ROBIE, C. *The Global Personality Inventory (GPI)*. In B. De Raad and M. Perugini, *Big Five Assessment*, p. 195–236, Alemanha, 2002.

SCHULTZ, D. *Theories Of Personality*. 4th ed. 1990.

SECOM. Pesquisa brasileira de mídia 2015: hábitos de consumo de mídia pela população brasileira, Secretaria de Comunicação Social, Brasília, 2015. ISBN: 978-85-85142-60-5. Disponível em: < <http://www.secom.gov.br/atuacao/pesquisa/>>.

SEIDMAN, G. *Self-presentation and belonging on Facebook: How personality influences social media use and motivations*. In: *Personality and Individual Differences*, v. 54, n. 3, p. 402-407, 2013.

SENTILEX-PT. SentiLex-PT. disponível em: <http://dmir.inesc-id.pt/project/SentiLex-PT_02>. Acesso em 25 abr. 2015.

SENTISTRENGTH. *SentiStrength*. Disponível em: < <http://sentistrength.wlv.ac.uk/>> . Acesso em: 21 abr. 2015.

SENTIWORDNET. *SentiWordNet*. Disponível em: < <http://sentiwordnet.isti.cnr.it/>> . Acesso em: 21 abr. 2015.

SHEVADE, S. K.; KEERTHI, S. S.; BHATTACHARYYA, C.; MURTHY, K. R. K. *Improvements to the SMO algorithm for SVM regression*. *IEEE Transactions on Neural Networks*, v. 11, n.5, p. 1188-1193, 2000.

SILVA, I. B.; NAKANO, T. D. C. Modelo dos cinco grandes fatores da personalidade: análise de pesquisas. *Avaliação psicológica*, p. 51-62, 2011.

SILVA, M. J.; CARVALHO, P.; SARMENTO, L. *Building a sentiment lexicon for social judgement mining*. In: *Computational Processing of the Portuguese Language*, Springer Heidelberg, p. 218-228, 2012.

SILVEIRA, B. A.; MURAMATSU, T. Y.; REVOREDO, K. C. Análise do perfil de uma comunidade científica através de mineração de texto, Dissertação de Mestrado, UFRJ, 2011.

SOUZA, C. V. R.; PRIMI, R.; MIGUEL, F. K. Validade do Teste Wartegg: correlação com 16PF, BPR-5 e desempenho profissional. *Avaliação psicológica*, v. 6, n. 1, p. 39-49, 2007.

SOUZA, D. A. NUNES, M. A. S. N. Aspectos teóricos e mensuração do construto psicológico personalidade desenvolvimento de plugin moodle para a formação de grupos de trabalho para uso na ead ufs, Relatório de pesquisa, 2011.

SOUZA, M.; VIEIRA, R. *Sentiment analysis on twitter data for portuguese language*. In: *Computational Processing of the Portuguese Language*, Springer Berlin Heidelberg, p. 241-247, 2012.

SOUZA, M.; VIEIRA, R.; Busetti, D.; CHISHMAN, R.; ALVES, I. M. *Construction of a Portuguese Opinion Lexicon from multiple resources*. In: *8th Brazilian Symposium in Information and Human Language Technology*, 2011.

SPEER, R.; HAVASI, C. *ConceptNet 5: A large semantic network for relational knowledge*. In *The People's Web Meets NLP*, Springer Heidelberg, Berlin, Germany, p. 161-176, 2013.

STRAPPARAVA, C.; VALITUTTI, A. *WordNet-Affect: An affective extension of WordNet*". In: *Proceedings of LREC*, p. 1083–1086, 2004.

SUMNER, C.; BYERS, A.; BOOCHEVER, R.; PARK, G. J. *Predicting dark triad personality traits from Twitter usage and a linguistic analysis of tweets*. In: *Machine Learning and Applications (ICMLA)*, v. 2, p. 386-393, 2012.

TAVARES, H. Personalidade, temperamento e caráter. In: *Fisiopatologia dos transtornos psiquiátricos*, Ed. Ateneu, São Paulo, p.191-205, 2006.

TENORIO, A. V. *Avaliação Qualitativa da Audiência de Televisão Baseada na Classificação de Sentimento de Usuários em Redes Sociais*. Dissertação de Mestrado, UFPE, 2014.

THAGARD, P. *Hot Thought: Mechanisms and Applications of Emotional Cognition*. MIT Press, Cambridge, MA, USA, 2006.

THELWALL, M. *Heart and soul: Sentiment strength detection in the social web with sentiment strength*. Cyberemotions, p. 1-14, 2013.

THELWALL, M.; BUCKLEY, K.; PALTOGLOU, G.; CAI, D.; KAPPAS, A. (2010). *Sentiment strength detection in short informal text*. In: Journal of the American Society for Information Science and Technology, p. 2544-2558, 2010.

TOMAÉL, M. I.; MARTELETO, R. M. Redes sociais: posição dos atores no fluxo da informação. Revista Eletr. de Bibliotecon, Florianópolis, 2006. Disponível em: <<https://periodicos.ufsc.br/index.php/eb/article/viewFile/342/387>>.

TOMLINSON, M. T.; HINOTE, D.; BRACEWELL, D. B. *Predicting Conscientiousness through Semantic Analysis of Facebook Posts*. In: Proceedings of WCPR13, Workshop on Computational Personality Recognition at ICWSM13 (7th International AAI Conference on Weblogs and Social Media), 2013.

TRAPPL, R.; PAYR, S.; PETTA, P. *Emotions in Humans and Artifacts*. MIT Press, Cambridge, MA, USA, 2003.

URQUIJO, S. Modelos circunplexos da personalidade, Contextos e questões da avaliação psicológica, São Paulo: Casa do Psicólogo, p. 31-49, 2001.

URSO, J. J. Stress e personalidade: overview e avaliação crítica de revisões sistemáticas sobre padrão comportamental tipo “a” e personalidade tipo “d” com desfechos coronarianos. Tese de Doutorado. Universidade de São Paulo. 2011.

UYSAL, I; GÜVENIR, H. A. *An overview of regression techniques for knowledge discovery*. In: The Knowledge Engineering Review, p. 319-340, 1999.

VALITUTTI, A.; STRAPPARAVA, C. *Interfacing Wordnet-affect with OCC model of emotions*. In: The Workshop Programme, p. 16, 2010.

WANG, Y. ; WITTEN, I. H. *Induction of model trees for predicting continuous classes*. In: Poster papers of the 9th European Conference on Machine Learning, 1997.

WEINERT, L. V. C. *Ontologias e técnicas de inteligência artificial aplicadas ao diagnóstico em fisioterapia neuropediátrica*, Tese de Doutorado (Universidade Tecnológica Federal do Paraná) - UTFPR, Curitiba, 2010.

WIEBE, J.; MIHALCEA, R. *Word Sense and Subjectivity*. In: Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics, p. 1065–1072, 2006.

WILSON, M. D. *The MRC Psycholinguistic Database: Machine Readable Dictionary, Version 2*. In: Behavioral Research Methods, Instruments and Computers, p. 6–11, 1988.

WITTEN, I. H., FRANK, E. *Data Mining: practical machine learning tools and techniques*. The Morgan Kaufmann Series in Data Management Systems, 2 ed. San Francisco: Elsevier, 2005.

WU, Y. *Disambiguating Dynamic Sentiment Ambiguous Adjectives*. n. August, p. 1191–1199, 2010.

YARKONI, T. *Personality in 100,000 words: A large-scale analysis of personality and word use among bloggers*. In: Journal of research in personality, p. 363-373, 2010.

YOUNG, L.; SOROKA, S. *Affective news: The automated coding of sentiment in political texts*. In: Political Communication, p. 205-231, 2012.

ZUO, X.; FENG, B.; YAO, Y.; ZHANG, T.; ZHANG, Q.; WANG, M.; ZUO, W. *A Weighted ML-KNN Model for Predicting Users' Personality Traits*. In: International Conference on Information Science and Computer Applications, 2013.

Apêndice A

Termo de consentimento para uso de informações pessoais

Neste apêndice é apresentado o documento de anuência concedido aos pesquisadores desse estudo, com a finalidade de coletar os dados nas redes sociais e acessar as respostas feitas no questionário de personalidade dos participantes do experimento.

TERMO DE CONSENTIMENTO LIVRE E ESCLARECIDO

Eu, nome: _____,
 nacionalidade: _____, idade: _____ anos, estado civil: _____,
 profissão: _____, endereço (rua, nº, bairro,
 cidade/UF e complemento): _____,
 _____, RG: _____,

estou sendo convidado a participar de um estudo denominado "Extração de Personalidade a partir de Textos em Português Brasileiro Utilizando Léxicos Linguísticos e Afetivos", cujo objetivo principal e justificativa são: desenvolver um método capaz de identificar cinco traços de personalidade contidas em textos escritos em português do Brasil. O Reconhecimento da Personalidade, consiste na classificação dos traços de personalidade de um indivíduo, expressada em diferentes mídias, dentre elas os textos. Nesta área, que ganhou impulso com a difusão da web, pesquisas vêm sendo desenvolvidas e grande parte delas visam criar métodos computacionais capazes de identificar traços de personalidade em textos. A maioria das pesquisas disponíveis atualmente é desenvolvida para a língua inglesa. Dessa forma, métodos que sejam capazes de identificar a personalidade em textos escritos em português do Brasil e que classifiquem a personalidade em traços (extroversão, neuroticismo, amabilidade, etc.) são uma contribuição relevante para a área.

A minha participação no referido estudo será no sentido de responder o inventário de personalidade NEO-IPIP (Neo-International Personality Item Pool). Esse inventário contém 120 questões em língua portuguesa, e para cada pergunta contém cinco opções de resposta (escala Likert: discordo totalmente, discordo parcialmente, nem discordo e nem concordo, concordo parcialmente e concordo totalmente). O inventário NEO-IPIP é maneira tradicional que psicólogos extraem os traços humanos e o instrumentos de mensuração de personalidade. A aplicação do inventário de personalidade tem o intuito de auxiliar a avaliação o método proposto, sem qualquer referência que possa identificar os voluntários da pesquisa. Estou ciente que os pesquisadores envolvidos na pesquisa terão acesso a informações publicadas do meu perfil das redes sociais: Facebook e Twitter, para fins de coletar textos que serão vinculados aos questionários de personalidade e utilizados para constituir uma base que será disponibilizada para a comunidade de pesquisadores da área. Tal base não contém identificação de autoria ou qualquer forma de reconhecimento que referencie a minha pessoa.

Fui alertado de que, da pesquisa a se realizar, posso esperar alguns benefícios, tais como: obter os traços de minha personalidade através de um inventário respaldado pela psicologia e com a construção de um método capaz de identificar automaticamente a personalidade de indivíduos por meio de textos escritos em português do Brasil, eu, assim como qualquer indivíduo da comunidade poderei usufruir desta pesquisa podendo aplicar o método em diversos contextos, como por exemplo, humor e expressão em figuras animadas

FORMIGADO LUIZ DE FERRAZ

FORMIGADO LUIZ DE FERRAZ

(avatars), extração de informações relacionadas a personalidade de indivíduos em textos postados em redes sociais, dentre outras aplicações.

Recebi, por outro lado, os esclarecimentos necessários sobre os possíveis desconfortos e riscos decorrentes do estudo, levando-se em conta que é uma pesquisa, e os resultados positivos ou negativos somente serão obtidos após a sua realização. Assim, estou ciente que o tempo despendido neste experimento poderá ser em vão, visto que a pesquisa pode não revelar resultados significativos.

Estou ciente de que minha privacidade será respeitada, ou seja, meu nome ou qualquer outro dado ou elemento que possa, de qualquer forma, me identificar, será mantido em sigilo.

Também fui informado de que posso me recusar a participar do estudo, ou retirar meu consentimento a qualquer momento, sem precisar justificar, e de, por desejar sair da pesquisa, não sofrer qualquer prejuízo à assistência que venho recebendo.

Os pesquisadores envolvidos com o referido projeto são prof. Dr. Fabrício Enembreck docente responsável pela pesquisa e membro permanente do Programa de Pós-Graduação em Informática da PUCPR e Aldo Marcelo Palm aluno de mestrado do Programa de Pós-Graduação em Informática da PUCPR, e com eles poderei manter contato pelo telefone (41) 3271-1669 ou (46) 8803-8653.

É assegurada a assistência durante toda pesquisa, bem como me é garantido o livre acesso a todas as informações e esclarecimentos adicionais sobre o estudo e suas consequências, enfim, tudo o que eu queira saber antes, durante e depois da minha participação.

Enfim, tendo sido orientado quanto ao teor de todo o aqui mencionado e compreendido a natureza e o objetivo do já referido estudo, manifesto meu livre consentimento em participar, estando totalmente ciente de que não há nenhum valor econômico, a receber ou a pagar, por minha participação.

Em caso de reclamação ou qualquer tipo de denúncia sobre este estudo devo ligar para o CEP PUCPR (41) 3271-2292 ou mandar um email para nep@pucpr.br

Curitiba, ____ de _____ de 2015.

Participante

Aldo Marcelo Palm

Prof. Dr. Fabrício Enembreck

Apêndice B

Resultados Suplementares de Experimentos

Este apêndice apresenta os resultados de vários experimentos complementares com o objetivo de avaliar o método proposto frente ao problema desta pesquisa.

As tabelas a seguir apresentam os resultados obtidos com os conjuntos TF-IDF normalizados, para cada traço do modelo *BigFive*.

Tabela 1: resultados obtidos com os conjuntos TF-IDF normalizados para o traço Extroversão.

Algoritmos	Extroversão								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
M5P	0.1734	0.113	0.1696	0.1389	0.2562	0.1365	0.1119	0.1776	0.2301
LinearRegression	0.174	-0.0008	-0.0289	-0.050	0.0404	0.0732	0.1745	0.2186	0.2689
SMOReg	0.2084	-0.0279	0.0652	-0.017	0.0433	0.0764	0.1822	0.1957	0.2639
SMOReg (Kernel = Puk)	0.096	0.0872	0.0901	0.0748	0.0769	0.0432	0.0439	0.024	0.0007
LWL	0.0172	0.0789	0.0185	-0.002	0.1097	0.1102	0.1815	0.1688	0.164
IBK (KNN=1)	0.021	0.018	0.0421	0.0001	0.1228	0.1076	0.1025	0.056	0.0798
IBK (KNN=3)	0.0161	0.043	0.0134	0.0019	0.011	0.0531	0.0708	0.1453	0.157

Tabela 2: resultados obtidos com os conjuntos TF-IDF normalizados para o traço Neuroticismo.

Algoritmos	Neuroticismo								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
M5P	0.0932	0.0365	-0.0084	0.0062	0.0233	0.1458	0.1419	0.0424	-0.0225
LinearRegression	0.0627	0.0465	0.0293	-0.026	0.0472	0.1078	0.0076	0.0109	-0.0344
SMOReg	0.0087	-0.0214	0.0578	0.0625	0.046	0.1073	0.0135	0.0074	-0.0372
SMOReg (Kernel = Puk)	0.0269	-0.0382	-0.0445	-0.074	-0.0166	-0.027	-0.1392	-0.143	-0.1517

Algoritmos	Neuroticismo								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
LWL	0.1622	0.137	0.1293	0.1149	0.0828	0.1038	0.0868	0.0368	0.0475
IBK (KNN=1)	0.0261	-0.0418	-0.0071	0.0152	-0.005	-0.008	0.0013	0.0246	-0.0623
IBK (KNN=3)	0.0975	0.0899	0.0679	0.0524	0.083	0.0654	0.0203	0.0205	-0.0446

Tabela 3: resultados obtidos com os conjuntos TF-IDF normalizados para o traço Realização.

Algoritmos	Realização								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
M5P	0.0543	-0.0614	0.0866	0.1011	0.1318	0.048	-0.0163	0.0918	-0.005
LinearRegression	0.0207	-0.0718	0.0377	-0.101	0.0377	0.0157	0.0059	0.0616	0.0922
SMOReg	-0.0213	0.0307	-0.0489	-0.022	0.0221	0.0035	-0.0048	0.044	0.0902
SMOReg (Kernel = Puk)	-0.0068	-0.0304	-0.0366	-0.035	-0.0194	-0.039	-0.0493	-0.054	-0.0676
LWL	0.0799	0.048	-0.0271	-0.000	0.0003	0.0201	0.0115	0.006	0.0406
IBK (KNN=1)	-0.0461	-0.0417	-0.0151	0.0176	0.1038	0.1125	0.0788	0.0513	0.0475
IBK (KNN=3)	-0.0147	-0.0667	-0.0811	-0.016	-0.0231	0.0564	0.0794	0.0411	0.0953

Tabela 4: resultados obtidos com os conjuntos TF-IDF normalizados para o traço Socialização.

Algoritmos	Socialização								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
M5P	-0.0276	-0.0602	0.0139	-0.043	0.0185	-0.018	0.0327	0.0652	0.0076
LinearRegression	0.1126	-0.0222	0.013	0.0676	0.0674	-0.012	0.0033	0.0344	0.0322
SMOReg	-0.0356	-0.0734	0.0829	0.0493	0.0535	-0.020	-0.0137	0.0291	0.0237
SMOReg (Kernel = Puk)	0.0346	0.0344	0.0473	0.0537	0.0561	0.0573	0.0459	0.0384	0.0265
LWL	-0.0705	0.0237	0.0382	0.0227	0.1232	0.0588	0.0919	0.1091	0.0854
IBK (KNN=1)	-0.0395	-0.083	-0.0167	0.0266	0.0262	-0.001	-0.0509	0.0324	0.077
IBK (KNN=3)	-0.1034	-0.0403	-0.0487	-0.013	-0.0112	0.0114	0.0927	-0.033	0.0967

Tabela 5: resultados obtidos com os conjuntos TF-IDF normalizados para o traço Abertura.

Algoritmos	Abertura								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
M5P	0.0595	0.0363	0.0155	0.1226	0.0832	0.0674	0.0983	0.0877	0.148
LinearRegression	0.1119	-0.0225	-0.0006	-0.028	0.0183	0.0499	0.0989	0.1771	0.1766
SMOReg	0.0903	-0.0336	-0.0824	0.0288	0.0183	0.0615	0.1035	0.1715	0.1726
SMOReg (Kernel = Puk)	0.0544	0.1085	0.1259	0.1298	0.1244	0.1174	0.1122	0.1033	0.0886

Algoritmos	Abertura								
	Número de Termos								
	68	150	200	250	300	500	750	1000	1500
LWL	0.0347	-0.013	0.0147	0.0824	0.0233	0.066	0.0162	0.015	0.0366
IBK (KNN=1)	-0.0882	0.0996	0.0863	0.1859	0.1438	0.0594	0.1183	0.1447	0.1723
IBK (KNN=3)	-0.0146	0.04	0.0617	0.0688	0.0888	0.1253	0.1495	0.1433	0.1509

As próximas tabelas apresentam os resultados obtidos da associação dos conjuntos TF-IDF com o léxico LIWC e os léxicos afetivos, sem aplicar a seleção de atributos na base de treinamento, para todos os traços de personalidade.

Tabela 6: Resultado do experimento combinando LIWC, léxicos afetivos e TF-IDF, sem aplicar a seleção de atributos na base de treinamento para o traço Extroversão.

Algoritmos	Extroversão								
	LIWC + Léxicos afetivos + Número de Termos TF-IDF								
	68	150	200	250	300	500	750	1000	1500
M5P	0.0526	0.0397	0.0762	0.045	0.095	0.1247	0.0915	0.1273	0.0526
LinearRegression	0.2033	0.1867	0.2018	0.1923	0.21	0.2089	0.2249	0.2388	0.2033
SMOReg	0.1884	0.1714	0.1794	0.1731	0.1965	0.1957	0.2073	0.2204	0.1884
SMOReg (Kernel = Puk)	0.0117	0.0074	0.0047	0.0005	-0.0018	-0.017	-0.0242	-0.035	0.0117
LWL	0.0924	0.1172	0.0939	0.0957	0.1595	0.153	0.1329	0.1393	0.0924
IBK (KNN=1)	0.0917	0.1065	0.1164	0.1462	0.1573	0.1026	0.1265	0.1433	0.0917
IBK (KNN=3)	0.0957	0.0938	0.0821	0.0992	0.1089	0.1094	0.082	0.1836	0.0957

Tabela 7: Resultado do experimento combinando LIWC, léxicos afetivos e TF-IDF, sem aplicar a seleção de atributos na base de treinamento para o traço Neuroticismo.

Algoritmos	Neuroticismo								
	LIWC + Léxicos afetivos + Número de Termos TF-IDF								
	68	150	200	250	300	500	750	1000	1500
M5P	0.0358	0.0785	0.1353	0.1336	0.1501	0.1701	0.072	0.0864	0.0358
LinearRegression	-0.0929	-0.0988	-0.0845	-0.080	-0.0595	-0.074	-0.0828	-0.068	-0.0929
SMOReg	-0.1037	-0.1064	-0.0899	-0.087	-0.063	-0.076	-0.0882	-0.073	-0.1037
SMOReg (Kernel = Puk)	-0.1617	-0.1579	-0.158	-0.155	-0.1545	-0.150	-0.1496	-0.148	-0.1617
LWL	0.1267	0.1311	0.1285	0.1215	0.1142	0.07	0.0749	0.0853	0.1267
IBK (KNN=1)	0.0115	0.0175	0.0256	0.0263	0.0531	-0.002	-0.0056	-0.012	0.0115
IBK (KNN=3)	-0.0412	-0.0371	-0.0704	-0.034	-0.0446	-0.041	0.0102	0.0065	-0.0412

Tabela 8: Resultado do experimento combinando LIWC, léxicos afetivos e TF-IDF, sem aplicar a seleção de atributos na base de treinamento para o traço Realização.

Algoritmos	Realização								
	LIWC + Léxicos afetivos + Número de Termos TF-IDF								
	68	150	200	250	300	500	750	1000	1500
M5P	0.1273	0.1387	0.1432	0.0906	0.0397	0.0457	0.0258	0.0431	0.1273
LinearRegression	0.0467	0.0483	0.0241	0.0356	0.0353	0.0225	0.0164	0.0457	0.0467
SMOReg	0.0492	0.0394	0.0199	0.0383	0.0339	0.0149	0.0008	0.0326	0.0492
SMOReg (Kernel = Puk)	-0.0613	-0.0633	-0.0623	-0.065	-0.0648	-0.072	-0.077	-0.078	-0.0613
LWL	0.0228	0.0206	0.0071	0.0071	0.0094	-0.003	0.007	0.0269	0.0228
IBK (KNN=1)	0.0471	0.0904	0.0813	0.0817	0.0841	0.062	0.0622	0.0604	0.0471
IBK (KNN=3)	-0.0433	-0.0278	-0.0376	-0.021	-0.0067	-0.017	-0.0112	0.04	-0.0433

Tabela 9: Resultado do experimento combinando LIWC, léxicos afetivos e TF-IDF, sem aplicar a seleção de atributos na base de treinamento para o traço Socialização.

Algoritmos	Socialização								
	LIWC + Léxicos afetivos + Número de Termos TF-IDF								
	68	150	200	250	300	500	750	1000	1500
M5P	-0.0354	0.0135	-0.011	0.0039	0.0033	-0.000	0.0199	0.0238	-0.0354
LinearRegression	-0.0499	-0.0711	-0.0607	-0.077	-0.0681	-0.070	-0.0473	-0.040	-0.0499
SMOReg	-0.0455	-0.0609	-0.0478	-0.057	-0.0564	-0.054	-0.0418	-0.032	-0.0455
SMOReg (Kernel = Puk)	0.0172	0.0138	0.0119	0.0107	0.0105	0.0089	0.0032	-0.001	0.0172
LWL	-0.0771	-0.0544	-0.0265	-0.054	0.0493	0.0581	0.0622	0.0545	-0.0771
IBK (KNN=1)	0.0902	0.0483	0.0482	0.0215	0.0046	0.0037	0.0205	0.0376	0.0902
IBK (KNN=3)	0.0081	0.0072	-0.0142	0.0102	-0.0115	0.0589	0.0999	0.0615	0.0081

Tabela 10: Resultado do experimento combinando LIWC, léxicos afetivos e TF-IDF, sem aplicar a seleção de atributos na base de treinamento para o traço Abertura.

Algoritmos	Abertura								
	LIWC + Léxicos afetivos + Número de Termos TF-IDF								
	68	150	200	250	300	500	750	1000	1500
M5P	0.0393	0.0407	0.0333	0.071	0.0575	0.0386	-0.0042	0.0156	0.0393
LinearRegression	0.0658	0.0634	0.067	0.0681	0.0761	0.066	0.0902	0.1371	0.0658
SMOReg	0.0716	0.0744	0.0801	0.0892	0.0986	0.0939	0.1064	0.1405	0.0716
SMOReg (Kernel = Puk)	0.0713	0.0691	0.0673	0.0631	0.0598	0.0497	0.0424	0.0364	0.0713
LWL	0.0689	0.0881	0.068	0.0976	0.099	0.0927	0.0932	0.0724	0.0689
IBK (KNN=1)	0.029	0.1142	0.142	0.1333	0.1287	0.1333	0.0835	0.0563	0.029
IBK (KNN=3)	0.0707	0.0729	0.085	0.1179	0.1039	0.0963	0.1085	0.1028	0.0707

As próximas tabelas apresentam os resultados obtidos da associação dos conjuntos TF-IDF com o léxico LIWC e os léxicos afetivos, com a aplicação de seleção de atributos na base de treinamento, para todos os traços de personalidade.

Tabela 11: Resultado do experimento combinando LIWC, léxicos afetivos e TF-IDF, com aplicação de seleção de atributos na base de treinamento para o traço Extroversão.

Algoritmos	Extroversão								
	LIWC + Léxicos afetivos + Número de Termos TF-IDF								
	68	150	200	250	300	500	750	1000	1500
M5P	0.6256	0.56	0.5997	0.5997	0.5829	0.5933	0.6079	0.653	0.6256
LinearRegression	0.706	0.6626	0.7661	0.8182	0.78	0.7489	0.8226	0.7649	0.706
SMOReg	0.4605	0.4672	0.5633	0.5633	0.6321	0.6321	0.6037	0.6104	0.4605
SMOReg (Kernel = Puk)	0.4561	0.4561	0.4561	0.4561	0.6208	0.5664	0.3876	0.6735	0.4561
LWL	0.3126	0.3126	0.3126	0.3126	0.3126	0.3126	0.2811	0.3061	0.3126
IBK (KNN=1)	0.1907	0.1907	0.1907	0.1907	0.1907	0.1907	0.2063	0.2063	0.1907
IBK (KNN=3)	0.2405	0.2405	0.2405	0.2405	0.2405	0.2369	0.2919	0.459	0.2405

Tabela 12: Resultado do experimento combinando LIWC, léxicos afetivos e TF-IDF com aplicação de seleção de atributos na base de treinamento para o traço Neuroticismo.

Algoritmos	Neuroticismo								
	LIWC + Léxicos afetivos + Número de Termos TF-IDF								
	68	150	200	250	300	500	750	1000	1500
M5P	0.5512	0.547	0.5531	0.5531	0.575	0.5579	0.5728	0.6107	0.5512
LinearRegression	0.7759	0.7229	0.7229	0.6402	0.7275	0.7553	0.8397	0.8476	0.7759
SMOReg	0.4511	0.4511	0.452	0.452	0.385	0.5165	0.3648	0.4713	0.4511
SMOReg (Kernel = Puk)	0.511	0.511	0.5111	0.5345	0.5345	0.5735	0.5763	0.5766	0.511
LWL	0.2731	0.2731	0.2731	0.2731	0.2731	0.2873	0.2873	0.1646	0.2731
IBK (KNN=1)	0.1685	0.1685	0.1685	0.1685	0.1685	0.1685	0.2524	0.3017	0.1685
IBK (KNN=3)	0.2764	0.2764	0.2764	0.2764	0.2764	0.2764	0.213	0.213	0.2764

Tabela 13: Resultado do experimento combinando LIWC, léxicos afetivos e TF-IDF, com aplicação de seleção de atributos na base de treinamento para o traço Realização.

Algoritmos	Realização								
	LIWC + Léxicos afetivos + Número de Termos TF-IDF								
	68	150	200	250	300	500	750	1000	1500
M5P	0.536	0.536	0.536	0.4958	0.4958	0.5589	0.6057	0.5752	0.536
LinearRegression	0.7638	0.7724	0.7415	0.6593	0.7999	0.8362	0.8145	0.7742	0.7638
SMOReg	0.4656	0.4762	0.5849	0.5439	0.4989	0.4952	0.5816	0.5354	0.4656
SMOReg (Kernel = Puk)	0.4663	0.466	0.4532	0.4532	0.4753	0.4754	0.4754	0.5488	0.4663
LWL	0.3168	0.3168	0.3168	0.3168	0.3168	0.3132	0.2302	0.394	0.3168
IBK (KNN=1)	0.3025	0.3025	0.3025	0.3025	0.3025	0.3025	0.3025	0.3065	0.3025
IBK (KNN=3)	0.2832	0.2832	0.2832	0.2832	0.2832	0.2832	0.2832	0.417	0.2832

Tabela 14: Resultado do experimento combinando LIWC, léxicos afetivos e TF-IDF, com aplicação de seleção de atributos na base de treinamento para o traço Socialização.

Algoritmos	Socialização								
	LIWC + Léxicos afetivos + Número de Termos TF-IDF								
	68	150	200	250	300	500	750	1000	1500
M5P	0.5082	0.5082	0.5082	0.5054	0.4874	0.473	0.5478	0.6055	0.5082
LinearRegression	0.6101	0.6684	0.6634	0.6491	0.7789	0.7986	0.7209	0.7203	0.6101
SMOReg	0.5	0.5	0.5169	0.4335	0.5645	0.5662	0.5225	0.3917	0.5
SMOReg (Kernel = Puk)	0.4422	0.4422	0.4422	0.4422	0.575	0.5505	0.5217	0.5933	0.4422
LWL	0.2243	0.2243	-0.0963	0.2243	0.2243	0.3372	0.3308	0.3419	0.2243
IBK (KNN=1)	0.1934	0.1934	0.1934	0.1934	0.1934	0.1934	0.1934	0.1934	0.1934
IBK (KNN=3)	0.2057	0.2057	0.2057	0.2057	0.2704	0.4302	0.4302	0.3261	0.2057

Tabela 15: Resultado do experimento combinando LIWC, léxicos afetivos e TF-IDF, com aplicação de seleção de atributos na base de treinamento para o traço Abertura.

Algoritmos	Abertura								
	LIWC + Léxicos afetivos + Número de Termos TF-IDF								
	68	150	200	250	300	500	750	1000	1500
M5P	0.5503	0.5503	0.5319	0.5319	0.5596	0.6142	0.6239	0.5924	0.5503
LinearRegression	0.7827	0.7635	0.7569	0.7569	0.7569	0.6618	0.8312	0.8967	0.7827
SMOReg	0.48	0.4662	0.4665	0.4665	0.4978	0.4658	0.4714	0.5172	0.48
SMOReg (Kernel = Puk)	0.4655	0.4655	0.4655	0.4623	0.4623	0.6072	0.5157	0.6354	0.4655
LWL	0.2745	0.2745	-0.1611	0.3221	0.3221	0.3221	0.3188	0.3188	0.2745
IBK (KNN=1)	0.2928	0.2928	0.2928	0.2928	0.2928	0.2623	0.2522	0.3125	0.2928
IBK (KNN=3)	0.445	0.445	0.445	0.445	0.445	0.445	0.445	0.4405	0.445

Anexo 1

Questões do Inventário NEO-IPIP 120

Na Tabela A.1, são apresentados todas as questões do inventário NEO-IPIP 120 de Johnson (2014), traduzidas por (NUNES; TELES; DE SOUZA, 2013) para o português brasileiro, utilizado no experimento. O questionário contém 120 questões. A tabela é composta por: (1) número da questão; (2) dimensão do *BigFive*ⁱ; (3) faceta correspondente e (4) texto da questão. Conforme apresentado a seguir:

Tabela A.1: Questões do inventário NEO-IPIP 120. Adaptado de (JOHNSON, 2014)

Número da Questão	Dimensão <i>BigFive</i> ⁱ	Faceta	Texto da questão
1	N1	Ansiedade	Preocupo-me com as coisas
2	E1	Amigabilidade	Faço amigos facilmente
3	A1	Imaginação	Tenho uma imaginação vívida
4	S1	Confiança	Confio nos outros
5	R1	Auto-eficácia	Completo tarefas com sucesso
6	N2	Raiva	Fico com raiva facilmente
7	E2	Gregarismo	Adoro festas com muitas pessoas
8	A2	Interesses artísticos	Acredito na importância da arte
9	S2	Moralidade	Nunca sonegaria impostos
10	R2	Ordem	Gosto de ordem
11	N3	Depressão	Frequentemente me sinto triste
12	E3	Assertividade	Assumo o comando das situações

Continua na próxima página

<i>Continuação da página seguinte</i>			
Número da Questão	Dimensão <i>BigFive</i>	Faceta	Texto da questão
13	A3	Emotividade	Vivo minhas emoções intensamente
14	S3	Altruísmo	Faço as pessoas se sentirem bem vindas
15	R3	Senso de dever	Tento obedecer as regras
16	N4	Auto-percepção	Sou intimidado facilmente
17	E4	Nível de atividade	Estou sempre ocupado
18	A4	Senso aventureiro	Prefiro variedade à rotina
19	S4	Cooperação	Sou fácil de satisfazer
20	R4	Empenho	Vou direto ao objetivo
21	N5	Autodisciplina	Frequentemente como demasiadamente
22	E5	Procura excitação	Adoro adrenalina
23	O5	Intelecto	Gosto de solucionar problemas complexos
24	S5	Modéstia	Detesto ser o centro das atenções
25	R5	Autodisciplina	Faço minhas tarefas o mais rápido possível
26	N6	Prudência	Entro em pânico com facilidade
27	E6	Bom humor	Irradio alegria
28	O6	Liberalismo	Tendo a votar em políticos de esquerda
29	S6	Compaixão	Tenho compaixão pelos desabrigados
30	R6	Prudência	Evito cometer erros
31	N1	Ansiedade	Tenho medo do pior
32	E1	Amigabilidade	Aproximo-me das pessoas com facilidade
33	A1	Imaginação	Curto altos vôos na minha imaginação
34	S1	Confiança	Acredito que os outros têm boas intenções
35	R1	Auto-eficácia	Sobressaio nas coisas que faço
36	N2	Raiva	Irrito-me facilmente
37	E2	Gregarismo	Converso com diversas pessoas em festas
38	A2	Interesses artísticos	Gosto de música
39	S2	Moralidade	Sigo as regras
40	R2	Ordem	Gosto de arrumar as coisas

Continua na próxima página

<i>Continuação da página seguinte</i>			
Número da Questão	Dimensão <i>BigFive</i>	Faceta	Texto da questão
41	N3	Depressão	Não gosto de mim mesmo
42	E3	Assertividade	Tento liderar os outros
43	A3	Emotividade	Sinto as emoções dos outros
44	S3	Altruísmo	Antecipo as necessidade dos outros
45	R3	Senso de dever	Mantenho as minhas promessas
46	N4	Auto-percepção	Tenho medo de fazer a coisa errada
47	E4	Nível de atividade	Estou sempre ativo
48	A4	Senso aventureiro	Gosto de conhecer lugares novos
49	S4	Cooperação	Não suporto confrontos
50	R4	Empenho	Trabalho duro
51	N5	Autodisciplina	Não sei porque faço algumas das coisas que faço
52	E5	Procura excitação	Busco aventura
53	O5	Intelecto	Adoro ler coisas que me desafiam
54	S5	Modéstia	Não gosto de falar sobre mim mesmo
55	R5	Autodisciplina	Estou sempre preparado
56	N6	Prudência	Muitas vezes me sinto sobrecarregado
57	E6	Bom humor	Divirto-me bastante
58	O6	Liberalismo	Acredito que não existe verdade absoluta
59	S6	Compaixão	Sinto compaixão por aqueles menos abastados que eu
60	R6	Prudência	Escolho minhas palavras com cuidado
61	N1	Ansiedade	Tenho medo de muitas coisas
62	E1	Amigabilidade	Sinto-me à vontade perto das pessoas
63	A1	Imaginação	Amo sonhar acordado
64	S1	Confiança	Confio no que as pessoas falam
65	R1	Auto-eficácia	Lido com minhas tarefas tranquilamente
66	N2	Raiva	Aborreço-me facilmente
67	E2	Gregarismo	Gosto de fazer parte de um grupo
68	A2	Interesses artísticos	Vejo beleza em coisas que outros podem não notar

Continua na próxima página

<i>Continuação da página seguinte</i>			
Número da Questão	Dimensão <i>BigFive</i>	Faceta	Texto da questão
69	S2	Moralidade	Uso de bajulação para avançar
70	R2	Ordem	Quero que tudo esteja perfeito
71	N3	Depressão	Frequentemente me sinto um lixo
72	E3	Assertividade	Convenço pessoas a agirem
73	A3	Emotividade	Sou apaixonado por causas
74	S3	Altruísmo	Adoro ajudar o próximo
75	R3	Senso de dever	Pago minhas contas em dia
76	N4	Auto-percepção	Tenho dificuldade de me aproximar das pessoas
77	E4	Nível de atividade	Faço diversas coisas no meu tempo livre
78	A4	Senso aventureiro	Interesso-me por muitas coisas
79	S4	Cooperação	Odeio parecer muito controlador ou exigente
80	R4	Empenho	Transformo planos em ações
81	N5	Autodisciplina	Faço coisas de que me arrependo posteriormente
82	E5	Procura excitação	Adoro ação
83	O5	Intelecto	Tenho um vocabulário rico
84	S5	Modéstia	Considero-me uma pessoa comum
85	R5	Autodisciplina	Inicio meus trabalhos o mais rápido possível
86	N6	Prudência	Sinto que sou incapaz de lidar com as situações
87	E6	Bom humor	Expresso alegria com uma criança
88	O6	Liberalismo	Acredito que criminosos deveriam receber ajuda ao invés de punição
89	S6	Compaixão	Valorizo mais cooperação do que competição
90	R6	Prudência	Sigo no caminho que escolho
91	N1	Ansiedade	Estresso-me facilmente
92	E1	Amigabilidade	Ajo confortavelmente perto de outras pessoas
93	A1	Imaginação	Gosto de me perder dos meus pensamentos
94	S1	Confiança	Acredito que as pessoas são essencialmente boas
95	R1	Auto-eficácia	Sei da minha capacidade
96	N2	Raiva	Estou frequentemente de mau humor

Continua na próxima página

<i>Continuação da página seguinte</i>			
Número da Questão	Dimensão <i>BigFive</i>	Faceta	Texto da questão
97	E2	Gregarismo	Envolvo outras pessoas no que estou fazendo
98	A2	Interesses artísticos	Amo flores
99	S2	Moralidade	Uso outras pessoas para conseguir meus objetivos
100	R2	Ordem	Gosto de ordem e harmonia
101	N3	Depressão	Tenho uma opinião ruim sobre mim mesmo
102	E3	Assertividade	Procuro influenciar outros
103	A3	Emotividade	Gosto de analisar a mim mesmo e minha vida
104	S3	Altruísmo	Preocupo-me com os outros
105	R3	Senso de dever	Falo a verdade
106	N4	Auto-percepção	Tenho medo de chamar atenção
107	E4	Nível de atividade	Consigo fazer muitas coisas ao mesmo tempo
108	A4	Senso aventureiro	Gosto de iniciar coisas novas
109	S4	Cooperação	Tenho uma língua afiada
110	R4	Empenho	Mergulho de coração nas minhas tarefas
111	N5	Autodisciplina	Gosto de farras
112	E5	Procura excitação	Gosto de fazer parte de multidões barulhentas
113	O5	Intelecto	Consigo lidar com muitas informações
114	S5	Modéstia	Raramente conto vantagem
115	R5	Autodisciplina	Começo logo a trabalhar
116	N6	Prudência	Não consigo me decidir
117	E6	Bom humor	Estou sempre de bem com a vida
118	O6	Liberalismo	Acredito numa única religião verdadeira
119	S6	Compaixão	Sofro com as perdas dos outros
120	R6	Prudência	Faço coisas sem pensar

ⁱ Dimensões do *BigFive* (N = Neuroticismo; E = Extroversão; S = Socialização; R = Realização; A = Abertura à experiência)

Anexo 2

Parecer consubstanciado do Comitê de Ética em Pesquisa da PUCPR

Neste apêndice é apresentado o documento submetido ao Comitê de Ética em Pesquisa da PUCPR, com a finalidade de obter o Certificado de Apresentação para Apreciação Ética (CAAE).



PARECER CONSUBSTANCIADO DO CEP

DADOS DO PROJETO DE PESQUISA

Título da Pesquisa: EXTRAÇÃO DE PERSONALIDADE A PARTIR DE TEXTOS EM PORTUGUÊS BRASILEIRO UTILIZANDO LÉXICOS LINGUÍSTICOS E AFETIVOS

Pesquisador: Fabrício Enembreck

Área Temática:

Versão: 2

CAAE: 45978515.1.0000.0100

Instituição Proponente: Pontifícia Universidade Católica do Paraná

Patrocinador Principal: Financiamento Próprio

DADOS DO PARECER

Número do Parecer: 1.175.821

Data da Relatoria: 05/08/2015

Apresentação do Projeto:

Conforme apresentado: "Este trabalho tem como objetivo principal o desenvolvimento de um método para reconhecimento da personalidade humana através de textos publicados online, escritos em língua portuguesa, por meio da análise de texto."

Objetivo da Pesquisa:

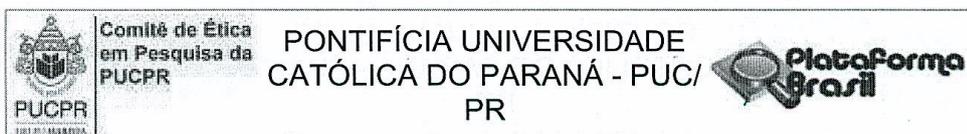
Conforme o projeto: "Objetivo Primário: Este trabalho tem como objetivo principal o desenvolvimento de um método para reconhecimento da personalidade humana através de textos publicados online, escritos em língua portuguesa, por meio da análise de texto."

Objetivo Secundário:

Os objetivos específicos incluem o levantamento bibliográfico sobre o reconhecimento de personalidade a partir de textos, o estudo de métricas de text-mining utilizando léxicos linguísticos e afetivos disponíveis em português, a mensuração dos traços de personalidade dos participantes aplicando inventários explícitos, a coleta dos textos públicos dos participantes, e implementação de um modelo que possibilite a mineração desses textos inferindo traços de personalidade por meio de pistas deixadas pelos autores, assim como uma avaliação empírica dos resultados, utilizando testes de hipótese não paramétricos."

Endereço: Rua Imaculada Conceição - 1155 - 3º andar
Bairro: Prado Velho **CEP:** 80.215-901
UF: PR **Município:** CURITIBA
Telefone: (41)3271-2103 **Fax:** (41)3271-2103 **E-mail:** nep@pucpr.br





Continuação do Parecer: 1.175.821

Avaliação dos Riscos e Benefícios:

Riscos:

Os dados extraídos das redes de relacionamento são públicos, portanto sua divulgação não apresenta risco aos usuários. Entretanto, a aplicação do Inventário NEO IPIP-120 produzirá dados sigilosos, que deverão ser de acesso restrito apenas aos pesquisadores.

Benefícios:

Um fator motivacional determinante à pesquisa é a expansão do Reconhecimento de Personalidade a partir de Textos para a língua portuguesa, considerando que existe carência de trabalhos que tratam da identificação de personalidade através de textos na área da computação para essa língua, tendo em vista que a grande maioria dos estudos se concentra no idioma inglês.

Comentários e Considerações sobre a Pesquisa:

A metodologia apresentada e a descrição dos elementos que compõem o projeto apresentam-se de forma satisfatória.

Considerações sobre os Termos de apresentação obrigatória:

Os termos apresentados estão conforme propostos pela legislação.

Recomendações:

Não há.

Conclusões ou Pendências e Lista de Inadequações:

A pendência relativa à postagem do inventário NEO IPIP 120 foi resolvida.

Situação do Parecer:

Aprovado

Necessita Apreciação da CONEP:

Não

Considerações Finais a critério do CEP:

Lembramos aos senhores pesquisadores que, no cumprimento da Resolução 466/12, o Comitê de Ética em Pesquisa (CEP) deverá receber relatórios anuais sobre o andamento do estudo, bem como a qualquer tempo e a critério do pesquisador nos casos de relevância, além do envio dos relatos de eventos adversos, para conhecimento deste Comitê.

Salientamos ainda, a necessidade de relatório completo ao final do estudo.

Eventuais modificações ou emendas ao protocolo devem ser apresentadas ao CEP-PUCPR de forma clara e sucinta, identificando a parte do protocolo a ser modificado e as suas justificativas.

Endereço: Rua Imaculada Conceição - 1155 - 3º andar
 Bairro: Prado Velho CEP: 80.215-901
 UF: PR Município: CURITIBA
 Telefone: (41)3271-2103 Fax: (41)3271-2103 E-mail: nep@pucpr.br



 <p>Comitê de Ética em Pesquisa da PUCPR</p>	<p>PONTIFÍCIA UNIVERSIDADE CATÓLICA DO PARANÁ - PUC/ PR</p>	
---	---	---

Continuação do Parecer: 1.175.821

CURITIBA, 07 de Agosto de 2015

Assinado por:
NAIM AKEL FILHO
(Coordenador)



Endereço: Rua Imaculada Conceição - 1155 - 3º andar
Bairro: Prado Velho CEP: 80.215-901
UF: PR Município: CURITIBA
Telefone: (41)3271-2103 Fax: (41)3271-2103 E-mail: nep@pucpr.br