

Análise de Sentimento de Tweets Relacionados aos Protestos que ocorreram no Brasil entre Junho e Agosto de 2013

Tiago C. de França¹, Jonice Oliveira,

Programa de Pós-Graduação em Informática da Universidade Federal do Rio de Janeiro (PPGI-UFRJ) – RJ – Brasil
tcruz.franca@ufrj.br, jonice@dcc.ufrj.br

Abstract. *The sentiment analysis of citizens is possible by using suitable techniques of analyzes applied to a massive database which is composed by messages provided by persons on Web. The goal of this paper is to analyze the opinion about protests that occurred in Brazil in 2013. For this, a database composed by tweets written in Brazilian Portuguese was used. This database was pre-processed for the corpus' creation. We observed that polarity (agreement or disagreement with the protests) of these messages and the final results have shown that the majority of messages are agreement ones.*

Resumo. *A análise de sentimento da população de um país é possível através da aplicação de técnicas adequadas sobre uma grande massa de dados formada por mensagens disponibilizadas pelas pessoas na Web. Este trabalho tem como objetivo analisar o sentimento acerca dos protestos que ocorreram no Brasil entre os meses de Junho e Agosto de 2013. Para tanto, foi criada uma base de tweets escritos em português brasileiro. Essa base foi pré-processada para criação do corpus de mensagens com menos ruídos. Esse corpus foi analisado para extração do sentimento presente nas mensagens. Observou-se a polaridade (apoio ou repúdio aos protestos) expressa nos tweets. Os dados foram analisados e o resultado final demonstrou que a maioria das mensagens apoiaram os protestos.*

1. Introdução

As mídias digitais da Web são fonte de uma vasta quantidade de informação disponibilizada em diferentes idiomas. Atualmente, as redes sociais na Web tem sido foco de diferentes tipos de estudo. Redes sociais *online* são redes formadas a partir da interação entre pessoas, grupos ou instituições motivadas por interesses ou objetivos comuns que se relacionam através de mídias digitais. É imensa a quantidade de usuários publicando informações de diferentes tipos e em diferentes idiomas nessas redes [Gonçalves et al. 2012; Nascimento et al. 2012].

Dentre as ferramentas que permitem a criação de redes sociais, está o Twitter¹, um *microblog* que permite que as pessoas divulguem qualquer tipo de informação quase em tempo real para todos aqueles ligados à sua rede. As publicações nessa plataforma são limitadas a um número pequeno de caracteres. Essa característica obriga que os usuários expressem sua opinião, sentimento ou qualquer informação através de mensagens curtas [Nascimento et al. 2012]. O Twitter possui mais de 200 milhões de usuários que geram aproximadamente 110 milhões de *tweets* por dia. Esses *tweets* possuem opiniões, informações pessoais ou sobre eventos em geral [Naaman e Boase 2010]. Por esse motivo, o Twitter tem sido visto como uma importante fonte

¹ <https://twitter.com>

para análise de opiniões e de sentimentos da população sobre eventos e acontecimentos [Li e Li 2011].

Neste trabalho foi realizada uma análise *tweets* relacionados aos protestos que se iniciaram no Brasil no mês de Junho de 2013. Nessa data o Brasil passou por momentos de protestos democráticos e muitas mensagens ligadas a esses acontecimentos foram publicadas no Twitter. Essas mensagens contêm informações relacionadas ao sentimento e a opinião da população acerca dos fatos.

O presente trabalho tem como objetivo analisar a polaridade expressa nos *tweets* relacionados aos protestos. Para tanto, uma base de *tweets* foi criada através da coleta de mensagens em períodos de notada relevância para o estudo, os quais ocorreram entre o mês de Junho e Agosto de 2013, com mais de 300 mil *tweets* (aproximadamente 1GB de conteúdo). Deseja-se verificar o apoio ou repúdio às manifestações através das opiniões presentes nas mensagens coletadas. A base de dados é constituída em sua maior parte por mensagens escritas no português brasileiro devido à forma como foram coletados os dados. Vale ressaltar que existem poucos trabalhos que tratam de escritos em português brasileiro.

A verificação da polaridade (apoio ou repúdio as manifestações) das mensagens da base *tweets* criada foi realizada por meio do emprego de algoritmos de aprendizagem estatística Naive Bayes. A escolha desses algoritmos se deu a partir da observação dos trabalhos na literatura que empregam essa mesma técnica de análises demonstram resultados satisfatórios [Lucca et al. 2013; Nascimento et al. 2012].

Este trabalho está organizado conforme segue. A Seção 2 contém um apanhado sobre os conceitos técnicos e teóricos necessários para realização do trabalho. A Seção 3 apresenta a descrição sobre a metodologia de desenvolvimento e as verificações realizadas. Na Seção 4 estão os resultados das análises realizadas. Finalmente, a Seção 5 apresenta as considerações finais e a descrição de possíveis trabalhos futuros.

2. Revisão de Literatura

O objetivo da análise de sentimentos (ou mineração de opinião) é definir, através da verificação dos termos que compõem um texto, se o documento analisado exprime uma opinião positiva, negativa ou neutra. Também é possível classificar qual sentimento está presente em um texto (por exemplo, raiva, felicidade, tristeza, etc.) [Karamibekr e Ghorbani 2012; Jambhulkar e Nirkhi 2014]. Esse tipo de análise é uma das áreas do processamento de linguagem natural que passou a ser mais investigada a partir dos anos 2000 [Liu 2012].

Muitos trabalhos de análise de sentimentos focam em textos escritos na língua inglesa pela grande quantidade de ferramentas, bases de conhecimento léxicas e ontológicas disponíveis para esse idioma. Porém, outros idiomas não possuem tantas bases ou mecanismos a disposição [Neri et al. 2012; Basile e Nissim 2013]. Este é o caso do português-brasileiro.

Embora haja uma quantidade razoável de autores brasileiros trabalhando com análises de sentimentos, muitos desses trabalhos, mesmo quando publicados na língua portuguesa, fazem análises de textos escritos em inglês. A principal razão pela pouca atuação de trabalhos que focam na análise de textos da língua portuguesa é a falta material de qualidade para realizar o pré-processamento e a análise, principalmente quando se trata de *tweets*, cujas mensagens são

ruidosas², com conteúdo informal, apresentando e muitas vezes erros gramaticais [Brew et al 2011; Neethu e Rajasree 2013].

Alguns trabalhos utilizam métodos de análise que se baseiam em *emoticons* (ícones ou sequência de caracteres que transmitem o estado emotivo da mensagem). Exemplos desses trabalhos são [Li and Li 2011; Zhang et al 2011; Hu et al. 2013]. Esse método de análise independe do idioma, porém não possuem grande abrangência, pois durante a análise apenas mensagens que possuem *emoticons* que estejam cadastrados na base de análise são considerados.

Silva et al. (2012) realizaram análise de sentimentos de textos utilizando regras gramaticais da língua inglesa e apontaram que sua abordagem pode ser utilizada para textos escritos em português, porém não foi demonstrado como isso pode ser feito. Além disso, o trabalho foi realizado sobre textos formais que seguem regras gramaticais e a norma culta da escrita o que não acontece com os *tweets*.

Freitas e Vieira (2013) utilizaram bases de ontologias para identificar a polaridade presente em mensagens em português relacionadas a publicações de usuários. Os autores conseguiram apenas bases ontológicas relacionadas a hotéis e filmes por não existirem bases disponíveis para o português. Tumitan e Becker (2013) apresentam um estudo de caso de mineração de opinião sobre comentários relacionados notícias sobre políticos em um jornal *online*. Para realizar a análise, os autores utilizaram uma base léxica do português de Portugal adaptado para o português brasileiro e para comentários ruidosos (com gírias, expressões próprias da Web e erros gramaticais). O trabalho não foi realizado sobre *tweets* e o dicionário léxico não pode ser reutilizado no presente trabalho. Os próprios autores apontam que em trabalhos futuros pretendem utilizar outras abordagens devido à dificuldade de se criar vocabulários específicos para o domínio do problema abordado.

Ferreira (2011) realizou a análise de documentos de empresas a fim de observar características presentes nesses documentos. O autor empregou alguns métodos de categorização, dentre eles o Naive Bayes o qual apresentou melhores resultados que as outras abordagens utilizadas. Apesar de não tratar de análise de sentimentos de *tweets*, esse autor demonstrou que Naive Bayes pode ser utilizado na análise conteúdo escrito em português.

O trabalho realizado por Varela (2012) buscou criar um sistema de classificação que seja independente da linguagem e realizou análises dos resultados através de bases de dados sobre cinema na língua espanhola e portuguesa. Os experimentos e análises realizados pelos autores demonstraram que o Naive Bayes se mostrou como técnica de análise mais adequada para textos pequenos. Shahheidari (2013) também utilizou Naive Bayes para o sentimento associado à *tweets* e constatou a eficácia desse método de análise. Nascimento et al. (2012) analisaram a reação positiva ou negativa da população após a divulgação de notícias pela mídia. Os autores coletaram dados de três tópicos pré-definidos de conteúdo em português, rotularam o conteúdo manualmente e realizaram a análise utilizando três métodos, entre eles o Naive Bayes o qual apresentou melhores resultados quando comparado aos outros dois métodos utilizados.

No que se refere a questões sociais, Bermingham e Conway (2009) utilizaram mensagens do YouTube para verificar a disponibilização e opinião sobre conteúdos ligados a grupos islâmicos radicais. Nagy et al. (2012) propuseram o uso de mecanismos de análise de sentimentos de *tweets* relacionados a crises ou desastres, enquanto Zhou et al. (2013)

² São mensagens que não seguem a norma culta de escrita ou regras gramaticais. Essas mensagens possuem gírias, palavras escritas erradas, expressões próprias da Web (“vc”, “lol”, “rsrsr” por exemplo), etc.

apresentaram um modelo de análise de sentimento de *tweets* para detectar interesses e opiniões da população sobre um evento social e testaram sua proposta analisando as eleições de 2010 da Austrália a fim de verificar a opinião da população sobre os candidatos. A IBM Research Brasil³ divulgou que está trabalhando sobre análises de *tweets* relacionados a jogos de futebol da seleção, porém não foram encontrados detalhes sobre a metodologia utilizada.

Dentre os trabalhos citados, pode-se verificar que nenhum deles buscou analisar sentimento relacionado aos protestos que ocorreram no Brasil em 2013. O objetivo do presente trabalho é analisar o sentimento da população brasileira que publicou mensagens no Twitter relacionados aos protestos de 2013 considerando o idioma oficial do país. A opinião pública pode influenciar as decisões políticas influenciando os tomadores de decisão assim como a opinião de consumidores influencia questões comerciais e fabricantes de produtos. Muitos dos trabalhos citados antes realizavam análise sobre conjuntos de dados com menos de 10 mil *tweets* ou de textos pequenos. A base analisada neste trabalho possui aproximadamente 1GB de tamanho antes do pré-processamento com mais de 300 mil *tweets*. Extrair informação de grande volume de dados é um desafio e traz novas oportunidades nesta área [Baumgarten 2013].

3. Metodologia

O objetivo deste trabalho é analisar a polaridade expressa nas mensagens publicadas no Twitter durante os protestos que ocorreram no Brasil nos meses de Junho, Julho e Agosto de 2013. A polaridade analisada busca investigar a quantidade de mensagens de apoio ou repúdio aos movimentos de protestos. Para realizar as análises, optou-se pelo uso de modelos estatísticos de aprendizagem Naive Bayes que é um sistema de classificação que independe de linguagem e que tem apresentado bons resultados na literatura [Lucca et al. 2013], [Nascimento et al. 2012], [Varela 2012].

O trabalho foi realizado seguindo duas etapas. A primeira etapa engloba a criação de uma base de *tweets* relacionados aos protestos e o pré-processamento dos dados coletados. A segunda etapa está relacionada ao método de análise utilizado. O algoritmo Naive Bayes precisa passar por um treinamento e sua acurácia pode ser verificada através de testes para comparação de resultados. Esta Seção descreve as ações realizadas desde a fase de pré-processamento até a fase da análise.

As funções desenvolvidas e utilizadas neste trabalho são apresentadas na Figura 1. Cada ação está enumerada em sequência de execução. As etapas de coleta, pré-processamentos, classificação manual, treino, teste e classificação automática dos *tweets* são realizadas em etapas distintas, porém cada função depende do resultado da função anterior com exceção da função de Coleta de dados.

3.1 - Coleta de dados e pré-processamento

A coleta de dados está relacionada com a recuperação das mensagens relacionadas aos protestos. No presente trabalho ela foi realizada com base em uma lista de *hashtags* relacionadas aos protestos. *Hashtags* são palavras chaves (palavras que designam um assunto abordado) antecedidas pelo caractere cerquilha (por exemplo, *#vemprarua*). A lista foi criada a partir da observação dos termos mais comuns utilizados para indexar os *tweets* os relacionando ao evento.

³ <http://www.research.ibm.com/index.shtml>

A lista de *hashtags* foi utilizada na consulta realizada através da API (do inglês, *Application Programming Interface*) disponibilizada pelo Twitter. O Twitter impõe algumas restrições dentre as quais estão o limite de acesso e a quantidade de *tweets* retornados por consulta. Os resultados das consultas foram armazenados em arquivos no forma JSON. Cada linha de JSON corresponde a uma ocorrência, um *tweet*. Foram realizadas cinco coletas de dados em períodos distintos entre os meses de Julho a Agosto. Os *tweets* retornados variam as datas de publicação do dia 30 de Junho ao dia 01 de Agosto do ano de 2013.

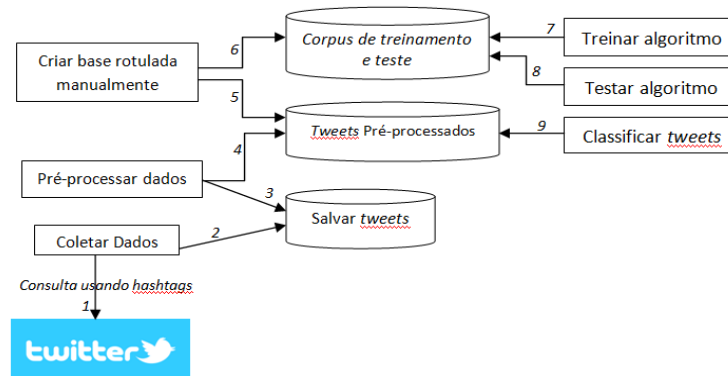


Figura 1 - Etapas da Análise de Polaridade dos tweets

De posse da base de dados, viu-se a necessidade de realizar o pré-processamento de todo conteúdo. No primeiro filtro foram retiradas as informações relevantes da base utilizada, pois alguns dados presentes em cada ocorrência de *tweet* (em cada linha de JSON) possuía um conjunto de informações irrelevantes para a análise desejada.

Em seguida, outro filtro foi aplicado para remover dos textos das mensagens menções a perfis do Twitter os quais são caracterizados por serem antecidos pelo caráter @ (*at*). Outro filtro removeu *retweets* (usuários publicando/compartilhando *tweets* de terceiros). Vale observar que os textos dos *retweets* são mantidos, pois foi considerado que ao *retweetar* uma mensagem, um usuário está concordando com o que ela expressa e demonstrando seu sentimento através da mesma. Caracteres de pontuação também foram removidos. Por último, foram removidas as URLs presentes nas mensagens, por ter sido considerado que URLs não representam nenhum sentimento. Todos os filtros possuem expressões regulares.

Depois dos filtros básicos, as mensagens passam por um procedimento de remoção de *stopwords*. Para isso, foi utilizado um conjunto de palavras categorizadas como *stopwords* do português brasileiro e, caso alguma das palavras estivesse presente na mensagem sendo analisada, tal palavra seria retirada do texto. A lista de *stopwords* utilizada é a mesma do mecanismo de busca Apache Solr⁴ e foi obtida nesse mecanismo.

O último pré-processamento é a aplicação do texto a uma operação de *stemming* que pode ser entendido como processo extração do radical de uma palavra. O método utilizado remove apenas o sufixo de uma palavra.

3.2 Classificação

O método de Análise adotado foi o Naive Bayes [Duda et al. 2000]. Naive Bayes é um método de aprendizagem probabilística baseado no teorema de Bayes. Este modelo assume que a

⁴ <https://lucene.apache.org/solr>

probabilidade de presença ou ausência de uma característica particular de uma classe não está relacionada com a presença ou ausência de qualquer outro recurso. Ou seja, o Naive Bayes é utilizado neste trabalho para realizar a classificação de textos baseada em sentimentos presente nas mensagens da base coletada.

O cálculo da probabilidade é simples e pode ser realizado da seguinte forma:

$$P[H | E] = (P[E | H] P[H]) / P[E]$$

onde H representa o evento observado, E a evidência, P[H] é a probabilidade do evento antes da evidência ser vista, e P [H | E] é a probabilidade de H dado o evento E (ou seja, a probabilidade de H após o evento E ser visto). P[E] representa o somatório da iteração de ambas as classes com a estimação da distribuição dos termos sobre as classes.

O algoritmo Naive Bayes tem sido demonstrado como sendo um método eficaz de análise. Para utilizá-lo é necessário apenas um pequeno grupo de treino durante a fase de aprendizado. Além disso, esse algoritmo é utilizado para calcular a probabilidade de uma mensagem ser categorizada dentro de uma classe pré-definida mesmo que os dados de entrada possuam ruídos [Ferreira 2012]. Os textos que são representados como *bag of words* para serem classificados com base em um modelo probabilístico. Isso significa que as posições exatas das palavras são ignoradas, e o classificador é montado com base no teorema de Bayes, assumindo independência entre as variáveis (palavras) [Jurafsky e Martin 2008]. Essas características justificam a adoção do algoritmo na presente análise.

No presente trabalho, o treino e os testes foram realizados a partir do uso de uma base rotulada por 3 humanos construída a partir de amostras de mensagens da base coletada. Ou seja, um conjunto de *tweets* pré-processados foi analisado por humanos e classificados como positivo (apoio aos protestos) ou negativo (repúdio aos protestos). Esse processo teve como objetivo cruzar os votos de cada um dos pesquisadores e eleger qual sentimento classificaria cada *tweet*. Ao final do processo, dois conjuntos ⁵(*corpus*) de 100 *tweets* cada foi criado. Um conjunto continham apenas mensagens positivas e outras negativas. Conforme Nascimento et al. (2012), não foram considerados *tweets* neutros, já que não foi possível encontrar na literatura um consenso sobre quais seriam as características típicas de textos classificados desta forma.

No presente trabalho, 70% da base rotulada por humano (*corpus* de treino) foi utilizada no treino. Enquanto 30% dessa base (*corpus* de teste) foram utilizados para teste e verificação da acurácia do algoritmo dado o conjunto de mensagens da base formada.

3.3. Análise

As análises visam verificar a qualidade do classificador, ou seja, a eficácia do classificador adotado. Para tanto, depois da etapa de treino, o classificador foi avaliado com a base de teste. Foi verificada a acurácia, a variância, o desvio, o desvio padrão, precisão, *recall* e *Macro-Avaraged* [Pak e Paroubek 2010; Nascimento et al. 2012 e Tsoumakas et al. 2010].

4. Análise dos Resultados

Esta Seção descreve como foram realizadas as etapas de coleta de dados, pré-processamento das mensagens, e a realização da análise da polaridade presente nas mensagens. Os resultados e análises sobre os resultados também são apresentados nesta Seção.

⁵ http://nltk.github.com/nltk_data

4.1. Criação das Bases

Os *tweets* da base foram coletados utilizando as *hashtags* apresentadas na Listagem 1. A coleta foi realizada com um *crawler* desenvolvido em Python⁶. Depois da extração apenas dos campos de interesse para a análise, a remoção de *stopwords* e o *stemming* foram implementados usando a linguagem Python e o módulo NLTK (*Natural Language Toolkit*). A base rotulada foi construída após essa etapa.

Listagem 1 - *Hashtags* usadas na coleta

#acordabrasil# OR #vempruarua OR #ForaFifa OR #ogiganteacordou OR #anonymousbrazil OR #MPL OR #passelivre OR #pec37 OR #mudabrasil OR #ChangeBrazil OR #anonymousbrazil OR #protesto OR #foradilma OR #protestorj OR #protestabrasil OR #primaverabrasileira OR #forafeliciano OR #ocupa OR #copapraquem OR #protest OR #pec33 OR #pec99

4.2. Verificação do Resultado do Classificador Naive Bayes

Nesta Seção apresentamos os resultados obtidos após a fase de treino e teste realizados sobre a base rotulada. Após treinar o algoritmo fornecendo ao NLTK 70% dos *tweets* classificados por humanos e rotulados como positivos ou negativos, utilizou-se os 30% restante para analisar a acurácia.

A Tabela 1 mostra: média de acerto, variância, desvio padrão, precisão, *recall*, *Macro-Average*, e *F-Socre* obtido após a execução dos testes. A Tabela 2 apresenta a matriz de confusão contendo as informações de *tweets* classificados corretamente ou de maneira equivocada. A base de comparação é formada por *tweets* rotuladas por humanos. As entradas da matriz de confusão são referidas como verdadeiros positivos (*true positives* - TP), falsos positivos (*false positives* - FP), falsos negativos (*false negatives* - FN) e verdadeiros negativos (*true negatives* - TN).

Tabela 1 – Acurácia (A), Variância (V), Desvio Padrão (DP), Precisão (P%), Recall (R%), Macro-Averaged (Ma-A) e F-score (F%) do classificador para as categorias testadas.

	A(%)	V	DP	P%	R%	Ma-A	F%
<i>Corpus</i> de Teste Positivo	90%	0.0325	0.1803	79%	87%	1.18	83%
<i>Corpus</i> de Teste Negativo	72%	0.0325	0.1803	85%	77%	1.05	81%

Tabela 2 - Matriz de Confusão

Categoria	Observação Real de Positivos	Observação Real de Negativos
Predição esperada positivo	26	4
Predição esperada negativo	7	23

Os resultados obtidos são satisfatórios, pois se a capacidade humana de avaliação correta da subjetividade de um texto varia de 72% [Wiebe et al 2006] a 85% [Golden 2011]. Se esse resultado for considerado como objetivo a ser alcançado, os experimentos demonstram que o objetivo foi atingido o que demonstra que o método escolhido é eficaz na classificação de *tweets*.

⁶ <http://www.python.org.br/>

4.3. Análise da polaridade do *tweets* na janela de tempo de coleta

Após avaliar o método adotado, foi realizada a análise de todos os dados da base coletada. O resultado da análise está apresentado na Figura 2. Essa figura apresenta a porcentagem de mensagens por polaridade positiva ou negativa. A reta azul representa a quantidade de mensagens na base que foram rotuladas como mensagens de apoio aos protestos. A reta vermelha indica a quantidade de mensagens presentes na base rotuladas como mensagens de repúdio aos processos. O eixo horizontal apresenta as datas da coleta e os dias entre as datas são os intervalos de dias de mensagens coletadas. Por exemplo, a coleta realizada no dia 30 de Junho obteve mensagens publicadas desde o dia 23 desse mês até o dia que estava sendo realizada a coleta.

Observando A Figura 2 é possível perceber que existe a indicação de que, com o passar do tempo as mensagens de apoio diminuíram e as de repúdio aumentaram. Isso pode ter ocorrido porque no Brasil aumentou a ocorrência de situações de violência, muitos confrontos e situações de destruição de patrimônio público e privado. Esses acontecimentos podem ter influenciado na opinião da população. No dia 23 os protestos que iniciaram por causa do aumento dos preços das passagens foram intensificados por causa das críticas as primeiras mobilizações quando foi dito que os protestantes não representavam a classe menos favorecida economicamente. Dias antes, no dia 18 de Junho, o “Anonymous Brasil” publica vídeo em resposta aos questionamentos da falta de foco dos protestos. Nos novos protestos entre 23 e 30 de junho as reivindicações incluíram uma série de insatisfações. Nessa data, os confrontos se intensificaram. No dia 11 de Julho. Outro motivo que pode estar ligado a diminuição do apoio aos protestos é o fim da Copa das Confederações no dia 30 de Junho [FIFA 2013] e o início do evento da Jornada Mundial da Juventude de 23 a 28 de Julho [JMJ 2013].

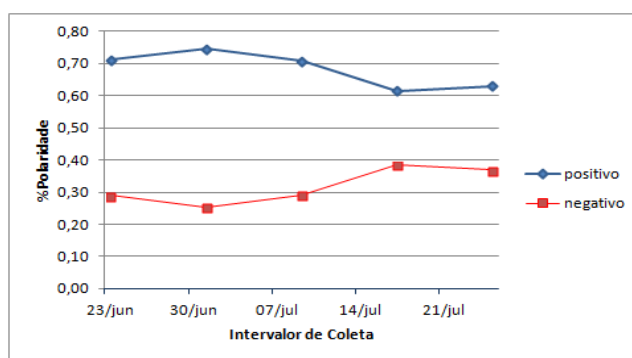


Figura 2 – Sentimentos das Mensagens Analisadas

4.4. Análise da polaridade por localização

Existem dois campos que podem ser utilizados para indicar a localização de um usuário. Um está associado ao usuário e ele preenche com texto livre. Esse campo pode ter preenchimento errado e é de difícil processamento por causa da falta de padrões. Por exemplo, um usuário pode escrever qual seu estado, outro usar apenas a abreviação para o estado, um terceiro usuário insere apenas sua cidade, etc. A outra forma de capturar a localização de um usuário é por meio da localização geográfica obtida do GPS do aparelho que ele utilizou para publicar a mensagem. Essa forma de localização possui um padrão. O Twitter utiliza as coordenadas obtidas do GSP do aparelho e identifica o país, a cidade e o estado do usuário. Neste trabalho, a segunda forma de localização foi utilizada para classificar os *tweets*.

Para agrupar os *tweets* por localização fornecida pelo GPS do dispositivo do usuário, é necessário que: o usuário possua um aparelho com GPS (um celular, por exemplo); e que ele permita que essa localização seja enviada para o Twitter. Por esse motivo, a maior parte das mensagens não possui localização. Aquelas que possuem foram agrupadas por Estados do Brasil. Outro grupo foi criado para contabilizar as mensagens publicadas fora do Brasil (chamado Internacional). O último grupo possui as mensagens que indicam que o país da pessoa que publicou é o Brasil, mas nenhum estado ou cidade estava presente no *tweet* impossibilitando assim em qual estado estava o usuário que publicou a mensagem. A Tabela 3 apresenta a quantidade de mensagens por grupo. A Tabela 4 apresenta a quantidade de *tweets* por Estado. O sentimento dos *tweets* classificados por Estado e os resultados foram agrupados por região do país.

Tabela 3 – Quantidade de Tweets por Grupo

Grupo	Sem Localização	Estados	Internacional	Outro
Quantidade de Mensagens	294755	6746	589	329

Tabela 4 – Quantidade (Qtde) de Tweets por Estado

Estado	Qtde	Estado	Qtde	Estado	Qtde
Acre	19	Maranhão	67	Rio de Janeiro	1261
Alagoas	46	Mato Grosso	22	Rio Grande do Norte	157
Amapá	38	Mato Grosso do Sul	37	Rio Grande do Sul	650
Amazonas	89	Minas Gerais	721	Rondônia	13
Bahia	234	Pará	252	Roraima	17
Ceará	196	Paraíba	145	Santa Catarina	207
Distrito Federal	313	Paraná	246	São Paulo	1542
Espírito Santo	177	Pernambuco	0	Sergipe	42
Goiás	161	Piauí	72	Tocantins	49

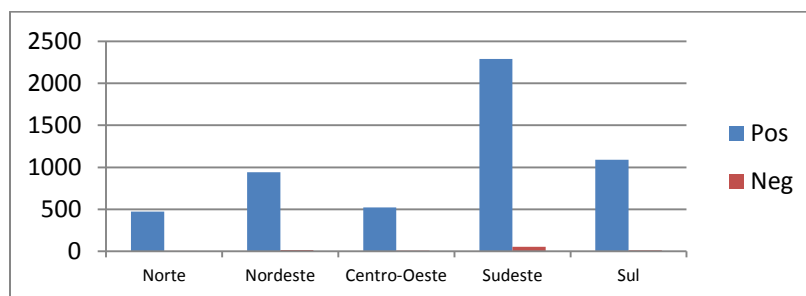


Figura 3 – Sentimentos das Mensagens por Região

6. Considerações finais

Este trabalho teve como objetivo analisar o sentimento de *tweets* relacionados aos protestos que ocorreram no Brasil em 2013. No trabalho foi utilizado o método de aprendizagem estatísticas Naive Bayes para analisar *tweets* em português brasileiro. Esse método independente do idioma e pode ser utilizado para analisar texto poluído como é característico dos *tweets*. Para isso, foi escolhido um classificador automático que depende de uma fase de treino supervisionado. Para treinar o classificador, *corpus* de *tweets* rotulados por humanos foram desenvolvidos sendo um positivo e outro negativo. Esse *corpus* foi construído utilizando *tweets*

coletados e pré-processados segundo a metodologia apresentada no trabalho. Uma parte de cada corpus (70%) foi utilizada para treino e a outra (30%) foi utilizada para testar o algoritmo e verificar a acurácia obtida.

Acredita-se que este trabalho se caracteriza como uma das poucas tentativas de realizar a associação entre *tweets* sentimentos (opiniões) relacionadas a eventos populares publicados no idioma português brasileiro para extrair o sentimento expresso pela população durante esses eventos.

Os resultados obtidos acerca do pré-processamento, treino e testes demonstrou resultados satisfatórios quando considerados aspectos de classificação realizada por humanos conforme demonstrado na literatura.

Após a verificação e aprovação dos resultados de treino e teste do algoritmo de análise de sentimentos adotado, foi realizada a verificação da polaridade de conjunto de mensagens coletadas durante o período observado. Verificou-se que grande parte das mensagens foram classificadas como mensagens de apoio (positivas) e que no final da coleta, houve uma pequena alteração na polaridade e mais mensagens de repúdio e menos de apoio foram publicadas. Porém o total de mensagens positivas ainda foi maior.

Esses resultados eram esperados por causa do que foi demonstrado por grande parte da população após observar as matérias acerca do assunto publicadas no período dos protestos. Ou seja, esperava-se que a quantidade de mensagens positivas fosse maior. Além disso, as coletas foram realizadas utilizando termos que estariam ligados diretamente aos protestos e não foram realizadas coletas considerando palavras-chave ou *hashtags* específicas de repúdio a mensagem (nenhum levantamento foi realizado para identificar se tais palavras ou *hashtags* existiram ou não).

Como trabalhos futuros se pretende criar um conjunto maior de treinos e testes empregando rotulação de grupos de pessoas a cega. Ou seja, cada pessoa rotula o *tweet* e o conjunto de mensagens rotuladas são comparadas para se chegar a um consenso. Além de aumentar o número de mensagens de treino e testes, também se pretende realizar análises mais detalhadas usando funções de distribuição acumulativa e teoria da decisão para verificar a robustez da polaridade das mensagens realizando possíveis ajustes de limiar. Por fim, pretende-se utilizar outros métodos de análise de sentimentos para comparar os resultados obtidos entre os métodos.

Referências

Basile, V.; Nissim, M. Sentiment analysis on Italian tweets, 2013. 4th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis 2013.

Baumgarten, M. Keyword-Based Sentiment Mining using Twitter, 2013. International Journal of Ambient Computing and Intelligence, 5 (2). pp. 56-69.

Birmingham, A.; Conway, M. Combining social network analysis and sentiment analysis to explore the potential for online radicalization. Disponível em: http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=5231878. Acessado em 10 de Março de 2014.

Brew, A., Greene, D., Archambault, D., and Cunningham, P. (2011). "Deriving Insights from National Happiness Indices.", 2011 IEEE 11th International Conference On Data Mining Workshops (ICDMW), pp. 53 –60.

Fantástico. Transporte e política são principais razões de manifestações, diz pesquisa. Disponível em: <http://g1.globo.com/fantastico/noticia/2013/06/transporte-e-politica-sao-principais-razoes-de-manifestacoes-diz-pesquisa.html>. Acessado em 10 de Fevereiro de 2014.

Ferreira, M. Classificação Hierárquica da Atividade Económica das Empresas a partir de Texto da Web, 2011. Disponível em <http://sigarra.up.pt/fep/pt/PUBLS_PESQUISA.FORMVIEW?p_id=13311>. Acessado em Setembro de 2013.

FIFA. Copa das Confederações, 2013. Disponível em: <http://pt.fifa.com/confederationscup/matches/>. Acessado em 30 de Março de 2014.

Freitas, L. A.; Vieira, R. Ontology based feature level opinion mining for portuguese reviews, 2013. Proceedings of the 22nd international conference on World Wide Web companion.

Gonçalves, P.; Dores, W.; e Benevenuto, F. PANAS-t: Uma Escala Psicometrica para Medição de Sentimentos no Twitter, 2012. Disponível em <http://www.imago.ufpr.br/csbc2012/anais_csbc/eventos/brasnam/artigos/>. Acessado em Agosto de 2013.

Hu, X.; Tang, J.; Gao, H.; Liu, H. Unsupervised Sentiment Analysis with Emotional Signals, 2013. Disponível em <<http://www.public.asu.edu/~xiahu/papers/www13.pdf>>. Acessado em Agosto de 2013.

Jambhulkar, P.; e Nirghi, S. A Survey Paper on Cross-Domain Sentiment Analysis. International Journal of Advanced Research in Computer and Communication Engineering Vol. 3, Issue 1, January 2014

JMJ. Jornada Mundial da Juventude, Rio 2013. Disponível em: <http://www.rio2013.com/>. Acessado em 30 de Março de 2014.

Jurafsky, D., and Martin, J. H. (2009). Speech and language processing: an introduction to natural language processing, computational linguistics, and speech recognition, Prentice Hall, 2nd Edition.

Karamibekr, M.; e Ghorbani, A. A. Sentiment analysis of social issues, 2012. Social Informatics (SocialInformatics).

Li, Y.-M., and Li, T.-Y. (2011). “Deriving Marketing Intelligence over Microblogs.”, Proceedings of 44th Hawaii International Conference On System Sciences (HICSS), pp. 1 –10.

Liu, B. Sentiment Analysis and Opinion Mining, Morgan & Claypool Publishers, Maio de 2012. Synthesis Lectures on Human Language Technologies.

Lucca, G.; Pereira, I. A.; Prisco, A.; Borges, E. N. Uma implementação do algoritmo Naïve Bayes para classificação de texto, 2013. Disponível em <<http://www.lbd.dcc.ufmg.br/colecoes/erbd/2013/0019.pdf>>. Acessado em Setembro de 2013.

Naaman, C.-H. L. Mor., and Boase, J. (2010). “Is it all About Me? User Content in Social Awareness Streams”, Proceedings of the 2010 ACM Conference on Computer Supported Cooperative Work, 2010.

Nagy, A.; Valley, C.M.S.; Stamberger, J. Crowd sentiment detection during disasters and crises. Proceedings of the 9th International ISCRAM Conference – Vancouver, Canada, April 2012.

- Nascimento, P.; Aguas, R.; Lima, D.; Kong, X.; Osiek, B.; Xexéo, G.; e Souza, J. Análise de sentimentos de tweets com foco em notícias, 2012. Disponível em <http://www.imago.ufpr.br/csbc2012/anais_csbc/eventos/brasnam/artigos> Acessado em Setembro de 2013.
- Neethu, M. S.; e Rajasree, R. Sentiment analysis in twitter using machine learning techniques. 2013 Fourth International Conference on Computing, Communications and Networking Technologies (ICCCNT).
- Neri, F.; Aliprandi, C.; Capeci, F.; e Cuadros, M. Sentiment Analysis on Social Media, 2012 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining.
- Pak, A., and Paroubek, P. (2010). “Twitter as a corpus for sentiment analysis and opinion mining.”, Proceedings of the 7th Conference on International Language Resources and Evaluation (LREC’10).
- Protestos no Brasil em 2013. Manifestações no Brasil em 2013. Disponível em: http://pt.wikipedia.org/wiki/Manifesta%C3%A7%C3%B5es_no_Brasil_em_2013. Acessado em 30 de Março de 2014.
- Shahheidari, S. Twitter Sentiment Mining: A Multi Domain Analysis. Complex, Intelligent, and Software Intensive Systems (CISIS), 2013 Seventh International Conference.
- Silva, N. R.; Lima, D.; e Barros, F. SAPair: Um Processo de Análise de Sentimento no Nível de Característica, 2012. Disponível em <<http://www.ppgia.pucpr.br/~enia/anais/wti/artigos.html>> Acessado em Setembro de 2013.
- Tumitan, D.; Becker, K. Tracking sentiment evolution on user-generated content: A case study in the brazilian political scene. Disponível em: <http://sbbd2013.cin.ufpe.br/Proceedings/application/06.html>. Acessado em 20 de Fevereiro de 2014.
- Varela, P.N. B. Sentiment Analysis, 2012. Disponível em <<http://goo.gl/PN3Xwg>>. Acessado em Setembro de 2013.
- Wiebe, J., Wilson, T., Cardie, C. (2006). “Annotating Expressions of Opinions and Emotions in Language”, Language Resources and Evaluation, v. 39, n. 2-3, pp. 165 –210.
- Wiebe, J., Wilson, T., Cardie, C. (2006). “Annotating Expressions of Opinions and Emotions in Language”, Language Resources and Evaluation, v. 39, n. 2-3, pp. 165 –210.
- Zhang, K., Cheng, Y., Xie, Y., Honbo, D., Agrawal, A., Palsetia, D., Lee, K., Liao, W., Choudhary, A. (2011). “SES: Sentiment Elicitation System for Social Media Data.”, Proceedings of 11th International Conference on Data Mining Workshops (ICDMW), pp. 129 – 136.
- Zhou, X.; Tao, X.; Yong, J.; Yang, Z. Sentiment analysis on tweets for social events. Computer Supported Cooperative Work in Design (CSCWD), 2013 IEEE 17th International Conference.
- Duda, R. O.; Hart, P. E.; e Stork, D. G. Patter Classification, Second Edtion, 2000. Editora TextBook.
- Tsoumakas, G.; Katakis, I.; Vlahavas, I. P. (2010). Mining multi-label data. In O. Maimon, L. Rokach (Eds.) *Data Mining and Knowledge Discovery handbook*, (pp. 667-685). Heidelberg, Germany: Springer-Verlag, 2nd Ed.